

0465

DUDLEY KNOX LIBRARY  
NAVAL POSTGRADUATE SCHOOL  
MONTEREY, CALIFORNIA 93940

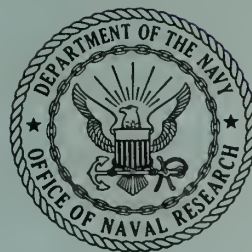
D.S.  
212310

# NAVAL RESEARCH LOGISTICS QUARTERLY

m

AM 11009  
CAL 11009

MARCH 1969  
VOL. 16, NO. 1



---

---

OFFICE OF NAVAL RESEARCH

NAVSO P-1278

# NAVAL RESEARCH LOGISTICS QUARTERLY

## EDITORS

Rear Admiral H. E. Eccles, USN (Retired)  
The George Washington University

O. Morgenstern  
Princeton University

F. D. Rigby  
Texas Technological College

D. M. Gilford  
U.S. Office of Education

S. M. Selig  
Managing Editor  
Office of Naval Research  
Washington, D.C. 20360

## ASSOCIATE EDITORS

R. Bellman, RAND Corporation  
J. C. Busby, Jr., Captain, SC, USN (Retired)  
W. W. Cooper, Carnegie Institute of Technology  
J. G. Dean, Captain, SC, USN  
G. Dyer, Vice Admiral, USN (Retired)  
P. L. Folsom, Captain, USN (Retired)  
M. A. Geisler, RAND Corporation  
A. J. Hoffman, International Business  
Machines Corporation  
H. P. Jones, Commander, SC, USN  
S. Karlin, Stanford University  
H. W. Kuhn, Princeton University  
J. Laderman, Office of Naval Research  
R. J. Lundegard, Office of Naval Research  
W. H. Marlow, The George Washington University  
R. E. McShane, Vice Admiral, USN (Retired)  
W. F. Millson, Captain, SC, USN  
H. D. Moore, Captain, SC, USN (Retired)

M. I. Rosenberg, Captain, USN (Retired)  
D. Rosenblatt, National Bureau of Standards  
J. V. Rosapepe, Commander, SC, USN (Retired)  
T. L. Saaty, U.S. Arms Control and  
Disarmament Agency  
E. K. Scofield, Captain, SC, USN  
M. W. Shelly, University of Kansas  
J. R. Simpson, Office of Naval Research  
J. S. Skoczylas, Colonel, USMC  
S. R. Smith, Naval Research Laboratory  
H. Solomon, The George Washington University  
I. Stakgold, Northwestern University  
E. D. Stanley, Jr., Rear Admiral, USN (Retired)  
C. Stein, Jr., Captain, SC, USN (Retired)  
R. M. Thrall, University of Michigan  
C. B. Tompkins, University of California  
J. F. Tynan, Commander, SC, USN (Retired)  
J. D. Wilkes, Department of Defense, OASD (ISA)

The Naval Research Logistics Quarterly is devoted to the dissemination of scientific information in logistics and will publish research and expository papers, including those in certain areas of mathematics, statistics, and economics, relevant to the over-all effort to improve the efficiency and effectiveness of logistics operations.

Information for Contributors is indicated on inside back cover.

The Naval Research Logistics Quarterly is published by the Office of Naval Research in the months of March, June, September, and December and can be purchased from the Superintendent of Documents, U.S. Government Printing Office, Washington, D.C. 20402. Subscription Price: \$1.50 a year in the U.S. and Canada, \$2.00 elsewhere; \$0.50 for a single copy. Requests for the purchase price of reprints of a particular article should be sent to the Superintendent of Documents, U.S. Government Printing Office, Washington, D.C. 20402.

The views and opinions expressed in this quarterly are those of the authors and not necessarily those of the Office of Naval Research.

Issuance of this periodical approved in accordance with Department of the Navy Publications and Printing Regulations, NAVSO P-35.

Permission has been granted to use the copyrighted material appearing in this publication.

# ON A NEW APPROACH TO THE ANALYSIS OF STATIONARY INVENTORY PROBLEMS\*

Oscar A. Z. Leneman

*Massachusetts Institute of Technology  
Lincoln Laboratory†  
Lexington, Massachusetts*

and

Frederick J. Beutler

*Computer, Information & Control Engineering Program  
The University of Michigan  
Ann Arbor, Michigan*

## ABSTRACT

The intent of this paper is to demonstrate that the theory of stationary point processes is a useful tool for the analysis of stationary inventory systems. In conventional inventory theory, the equilibrium distributions for a specified inventory policy are obtained, whenever possible, by recursive or limiting procedures, or both. A different and more direct approach, based on stationary point processes, is proposed here. The time instants at which stock delivery is effected are viewed as points of the stationary point process, which possesses uniform statistical properties on the entire real axis; hence the equilibrium statistics of the inventory process can be calculated directly. In order to best illustrate this approach, various examples are given, including some that constitute new results.

## INTRODUCTION

Once a class of inventory policies is specified, the study of the inventory model becomes a study of the associated stochastic process. Quite often one seeks to choose within this class the policy which minimizes the average long run costs and as a result one is interested in studying the stationary behavior of the inventory model. This type of study is, in general, difficult and the usual approach, utilizing recursive and limiting procedures [1,4,5] is often limited.

The intent of this paper is to show by examples and illustrations that such stationary state inventory studies are systematically facilitated by the use of some recent results from the theory of stationary point processes (hereafter abbreviated s.p.p.). An s.p.p. possesses uniform stationarity properties over the entire real axis, so that an inventory process based

---

\*This study was sponsored by Lincoln Laboratory, a center for research operated by the Massachusetts Institute of Technology with support from the U.S. Air Force and from the U.S. Advanced Research Projects Agency. Support also furnished by the National Aeronautics and Space Administration under Research Grant NsG-2-59. This article was originally drafted by O. A. Z. Leneman and revised for publication by F. J. Beutler during the period of leave in India of the first author.

†Operated with support from the U.S. Advanced Research Projects Agency.

on an s.p.p. will be in the same statistical state for all times  $t$ ,  $-\infty < t < +\infty$ . The statistics of the inventory process can thus be calculated directly; the final result does not involve  $t$ , and no asymptotic techniques are required.

For the reader unfamiliar with s.p.p. we state some of the s.p.p. definitions and properties, mentioning in particular those needed for the examples that follow. A detailed exposition of s.p.p. is found elsewhere [2,3], and the reader desiring proofs or further information should refer to [2,3]. A random point process  $\{t_n\}$  is a denumerable ordered family of random variables that may be regarded as points (occurrences) on the time axis. In particular, the points  $t_n$  can be related to inventory processes by taking these points to be the instants at which orders are placed for each successive inventory cycle.

An s.p.p. is a random point process meeting certain statistical uniformity requirements. For any finite set of intervals on the real axis, the joint probability distribution function of numbers of points on these intervals must be invariant under any translation of the set of intervals. This definition of stationarity is equivalent to any one of a number of other possible definitions. For instance, if  $L_k(t)$  is the time interval extending from time  $t$  to the  $k$ 'th point on the right of  $t$ , stationarity of the point process is equivalent to the following: for any  $n$ , the joint distribution of the set of random variables  $\{L_1(t), L_2(t), \dots, L_n(t)\}$  does not depend on  $t$ . It should be noted that the intervals  $X_n = t_{n+1} - t_n$  of an s.p.p.  $\{t_n\}$  need be neither independent nor identically distributed.

The second definition of an s.p.p. makes it possible to define a probability distribution function  $G_n(x) = P[L_n(t) \leq x]$ . The random variable  $L_{-n}(t)$ , defined as the time interval extending from the  $k$ 'th point to the left of  $t$  to time  $t$ , has the same distribution function  $G_n$  as  $L_n$ . It can be shown that the  $G_n$  are absolutely continuous, and that the derivatives  $g_n$  possess limits from the right. The  $G_n$  have a useful characterization for small argument. Suppose that  $\lim_{x \rightarrow 0} [G_2(x)/G_1(x)] = 0$ , that is, the probability of two points (or more) in a small interval is negligible when compared with the probability of only one. Such an assumption is certainly appropriate to inventory processes, where reordering frequency is subject to natural limitations. We then have  $\sum_{n=2}^{\infty} G_n(x) = o(x)$  and

$$(1) \quad G_1(x) = \beta x + o(x),$$

in which  $\beta$  represents the mean number of points per unit time. Because the sum of the  $G_n$ ,  $n \geq 2$ , is negligible by comparison,  $G_1(x) = P\{\text{one point falls in } (t, t+x]\}$  for small  $x$ .

Another distribution function that will often appear in the examples below is  $F_k(x) = P[L_k(t) \leq x | t = t_n \text{ for some } n]$ . When the intervals  $X_n$  are identically (but not necessarily independently) distributed,  $F_k$  is the probability distribution function for  $k$  successive intervals, that is

$$(2) \quad F_k(x) = P[X_1 + X_2 + \dots + X_k \leq x].$$

We note further that  $F_k$  obeys the relation

$$F_k(x) = 1 - \beta^{-1} \sum_{j=1}^k g_j(x) \quad \text{for } x \geq 0,$$



and that we can use  $F_1$  to show that the expectation

$$(3) \quad E(X_n) = \beta^{-1}.$$

In the examples that follow, we shall often assume that the stock depletion rate  $\lambda$ , the time delay  $\tau$  between ordering and delivery, and the amount  $M$  delivered are independent random variables. Each of these correspond to realistic situations; for instance, the demand rate  $\lambda$  is not precisely known in advance, and the delivered quantity  $M$  may represent the imprecisely realized output of an entire plant or farm, and may be subject to an unpredictable number of rejected units on receiving inspection or a random number of units damaged in shipment. In any case, considering  $\lambda$ ,  $\tau$ , and  $M$  to be random generalizes any analysis in which they are taken to be deterministic.

#### EXAMPLE I: A CONTINUOUS TIME INVENTORY POLICY

The following example will serve as an introduction. An amount of stock  $M$  is re-ordered with instantaneous delivery whenever the inventory on hand reaches zero. The stock on hand after delivery is depleted at a linear rate  $\lambda$  per unit time. The quantities  $M$  and  $\lambda$  are independent random variables selected from probability density functions  $f_M$  and  $f_\lambda$  for each inventory cycle. Assuming that the system has reached steady state, what is the probability density function of  $X(t)$ , the stock on hand at any given time  $t$ ?

#### SOLUTION

The stock on hand is a stationary random process  $X(t)$  as depicted in Figure 1. The instants of stock arrival at the random times  $t_n$  constitute a stationary point process. Successive  $t_n$  are separated by identically distributed intervals; if  $T$  is such an interval

$$(4) \quad T = \frac{M}{\lambda}.$$

The average number of stock arrivals per unit time is computed from (3) to be

$$(5) \quad \beta = \frac{1}{E(T)} = \frac{1}{E(M)E(1/\lambda)},$$

where  $E(\ )$  denotes the expectation of the quantity in parentheses.

Let us introduce the notation  $A_0 = \{\text{no stock arrival in } (t, t+y]\}$  and  $A_1 = \{\text{one stock arrival in } (t+y, t+y+dy]\}$ ; these two events will be mentioned several times in what follows. Now

$$(6) \quad P[\{x < X(t) \leq x + dx\} \cap \{y < L_1(t) \leq y + dy\}] = P[\{x < X(t) \leq x + dx\} \cap A_0 \cap A_1] \\ = P(A_1)P[\{x < X(t) \leq x + dx\} \cap A_0 | A_1].$$

The first probability on the right is obtained from (1). If  $A_1$  and  $A_0$  are true,  $X(t)$  must be given by  $\lambda y$ , and  $T > y$  from the argument on interval distributions preceding (2). Hence the probability (6) becomes successively

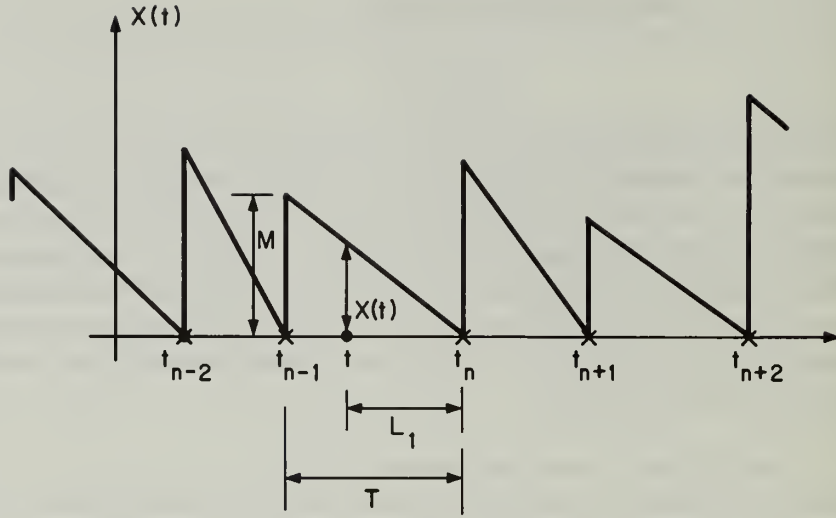


Figure 1. Inventory process with random linear depletion and instantaneous random reordering

$$\begin{aligned}
 (7) \quad \beta \, dy \, P[\{x < \lambda y \leq x + dx\} \cap \{T > y\}] &= \beta \, dy \, P\left[\left\{\frac{x}{y} < \lambda \leq \frac{x}{y} + \frac{dx}{y}\right\} \cap \left\{\frac{M}{\lambda} > y\right\}\right] \\
 &= \beta \, dy \, P\left[\frac{x}{y} < \lambda \leq \frac{x}{y} + \frac{dx}{y}\right] P[M > x] \\
 &= \beta \, dy \, f_{\lambda}\left(\frac{x}{y}\right) \frac{dx}{y} [1 - F_M(x)],
 \end{aligned}$$

in which  $F_M$  is the probability distribution function for  $M$ . In order to find the density of  $X(t)$ , one integrates [see the left side of (6)] on  $y$  for the marginal density. Thus

$$\begin{aligned}
 (8) \quad P[x < X(t) \leq x + dx] &= \beta \, dx [1 - F_M(x)] \int_0^{\infty} y^{-1} f_{\lambda}\left(\frac{x}{y}\right) dy \\
 &= \beta \, dx [1 - F_M(x)] \int_0^{\infty} u^{-1} f_{\lambda}(u) du \\
 &= \beta \, dx [1 - F_M(x)] E(\lambda^{-1})
 \end{aligned}$$

which yields the rather simple final result

$$(9) \quad f_X(x) = [1 - F_M(x)]/E(M).$$

The density of  $X(t)$  does not depend on the statistics of  $\lambda$ , because  $\lambda$  affects the rate of ordering and variations in that rate, but not the amplitude ordered or the fraction of that amplitude on hand at a given instant.

#### EXAMPLE II: A DISCRETE TIME (s,S) POLICY

This example, which is well known in the literature [1], will also serve as an illustration of our approach. Note that it is no longer necessary to show the existence of the steady state (see [1], pp. 234-237) when using s.p.p. techniques.

Whenever the stock level falls below  $s$ , ordering is immediately enacted to raise the level to  $S$  with immediate delivery; when the quantity in supply exceeds  $s$ , no ordering is done. A negative stock level is admissible, and can be viewed as an amount owed to consumption. Demands on the stock occur in discrete periods of duration  $T_0$ , the random demand sequence consisting of mutually independent identically distributed variates. The appearance of the stock level  $X(t)$  as a function of time  $t$  is then as shown in Figure 2. Our aim is again to determine the probability density  $f_X$  of the stock on hand  $X(t)$ .

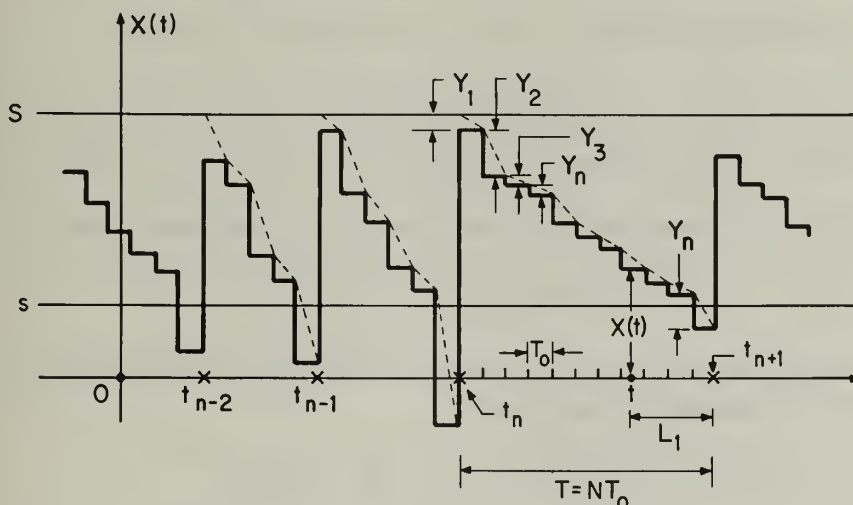


Figure 2. Inventory process with an (s,S) policy

#### SOLUTION

The  $t_n$  representing ordering times again constitute an s.p.p. with independent identically distributed intervals. Any time  $t$  will be in some such interval, and we take the length of that interval to be the random variable

$$(10) \quad T = NT_0 ;$$

here  $N$  is the number of periods of length  $T_0$  required to reduce the stock from  $S$  (at the ordering time immediately preceding  $t$ ) to  $s$  or less. Let  $Y_1, Y_2, \dots$ , be the random demands for this ordering interval. The total demand for the first  $n$  intervals is then defined to be  $S_n = Y_1 + Y_2 + \dots + Y_n$ . This means that the random integer  $N$  satisfies the inequalities

$$(11) \quad S_{N-1} < S - s \leq S_N.$$

As in the preceding example, we shall need the average number of stock arrivals per unit time  $\beta$ . From the relationship (3) between  $\beta$  and the interval length (10), we have

$$(12) \quad \beta = \frac{1}{E(T)} = \frac{1}{T_0 E(N)}.$$

It remains to calculate  $E(N) = \sum_{n=1}^{\infty} nP[N=n]$ . Now

$$(13) \quad P[N=n] = Q_{n-1}(S-s) - Q_n(S-s),$$

where  $Q_n$  is the probability distribution function of  $S_n$  with the convention that  $Q_0$  is the unit step function. Hence

$$(14) \quad E(N) = 1 + H(S-s),$$

where  $H(t)$  is the renewal function

$$(15) \quad H(t) = \sum_{n=1}^{\infty} Q_n(t).$$

From (12) and (14) the average number of stock arrivals per unit time is therefore

$$(16) \quad \beta = \{T_0[1 + H(S-s)]\}^{-1}.$$

The density function  $f_X$  of  $X(t)$  is computed separately for the two cases  $x \leq s$  and  $s < x \leq S$ . We start with the former. Consider

$$(17) \quad P[\{x < X(t) \leq x + dx\} \cap \{y < L_1(t) \leq y + dy\} \cap \{N=n\}] \\ = P[\{x < S - S_n \leq x + dx\} \cap \{S_{n-1} < S - s \leq S_n\} \cap A_0 \cap A_1].$$

The right-hand side was obtained by inserting the condition for  $L_1$  in terms of  $A_0$  and  $A_1$ , and by using (11) to rewrite the set  $\{N=n\}$ . We pursue the same argument as in (6) and (7) in Example I. More specifically, the probability (17) can be written

$$(18) \quad \beta dy P[\{S - x - dx < S_n \leq S - x\} \cap \{S_{n-1} \leq S - s\} \cap A_0 | A_1] \\ = \beta dy P[\{S - x - dx < S_n \leq S - x\} \cap \{S_n \leq S - s\}]$$

for  $0 < y \leq T_0$ . If  $n=1$ , the set  $\{S_n \leq S - s\}$  does not appear, and so

$$(19) \quad P[\{x < X(t) \leq x + dx\} \cap \{y < L_1(t) \leq y + dy\} \cap \{N=1\}] \\ = \beta dy P[S - x - dx < Y_1 \leq S - x] = \beta dy q_1(S - x) dx.$$



In (19),  $q_n$  is the probability density for the sum of any  $n$  of the  $Y_k$ , that is,  $q_n(x) = dQ_n(x)/dx$ . The marginal probability

$$(20) \quad P[\{x < X(t) \leq x + dx\} \cap \{N = 1\}] = \beta T_0 q_1(S - x) dx$$

is found from (19) by integrating on  $y$  from zero to  $T_0$ , where the upper limit is the largest value attainable for  $L_1(t)$  when  $X(t) \leq s$ .

Now for  $n=2, 3, \dots$ , the probability (17) is first developed according to (18). When the right side of (18) is expressed in terms of the probability densities  $q_n$  there results

$$(21) \quad \begin{aligned} P[\{x < X(t) \leq x + dx\} \cap \{y < L_1(t) \leq y + dy\} \cap \{N = n\}] \\ = \beta dy \int_{u \leq S-s} P[\{u < S_{n-1} \leq u + du\} \cap \{S - x - u < Y_n \leq S - x - u + du\}] \\ = \beta dy dx \int_0^{S-s} q_{n-1}(u) q_1(S - x - u) du. \end{aligned}$$

The right side of (21) has used the assumption of mutual independence and identical distributions for the demands  $Y_k$ . The reasoning leading to (20) yields for the marginal probability

$$(22) \quad P[\{x < X(t) \leq x + dx\} \cap \{N = n\}] = \beta T_0 dx \int_0^{S-s} q_{n-1}(u) q_1(S - x - u) du.$$

The formula for total probability

$$(23) \quad P[x < X(t) \leq x + dx] = \sum_{n=1}^{\infty} P[\{x < X(t) \leq x + dx\} \cap \{N = n\}]$$

can now be applied to the joint probabilities (19) and (22). Thus

$$(24) \quad f_X(x) = \beta T_0 \left\{ q_1(S - s) + \int_0^{S-s} q_1(S - x - u) \sum_{n=1}^{\infty} q_n(u) du \right\}.$$

The sum on the right is the derivative of the renewal function  $H(u)$ , and will be denoted by  $h(u)$ . The desired result is therefore

$$(25) \quad f_X(x) = \frac{1}{1 + H(S - s)} \left\{ q_1(S - s) + \int_0^{S-s} q_1(S - s - u) h(u) du \right\}$$

by a substitution for  $\beta$  from (16).

The second case,  $s < x \leq S$ , remains to be considered. For this case  $n \geq 2$ , and  $kT_0 < y \leq (k+1)T_0$  for some  $k=1, 2, \dots, n-1$ . Now

$$\begin{aligned}
 (26) \quad & P[\{x < X(t) \leq x + dx\} \cap \{y < L_1(t) \leq y + dy\} \cap \{N = n\}] \\
 &= P[\{x < S - S_{n-k} \leq x + dx\} \cap \{S_{n-1} < S - s \leq S_n\} \cap A_0 \cap A_1] \\
 &= \beta dy P[\{S - x - dx < S_{n-k} \leq S - x\} \cap \{S_{n-1} < S - s \leq S_n\} \cap \{T > y\}].
 \end{aligned}$$

Since  $\{T > y\}$  is already implied by the intersection of the other sets on the right of (26) by the choice of  $k$

$$\begin{aligned}
 (27) \quad & P[\{x < X(t) \leq x + dx\} \cap \{y < L_1(t) \leq y + dy\} \cap \{N = n\}] \\
 &= \beta dy q_{n-k}(S - x) dx [Q_{k-1}(x - s) - Q_k(x - s)].
 \end{aligned}$$

When the left side of (27) is summed over  $n$  (starting with  $n = k + 1$  because  $k + 1$  is the minimum value that can be assumed by  $N$ ) the total probability law yields

$$\begin{aligned}
 (28) \quad & P[\{x < X(t) \leq x + dx\} \cap \{y < L_1(t) \leq y + dy\}] \\
 &= \beta dy dx [Q_{k-1}(x - s) - Q_k(x - s)]h(S - x).
 \end{aligned}$$

As before,  $h(u) = \sum_{n=1}^{\infty} q_n(u)$ . Next, let us integrate (28) on  $y$  over the interval  $(kT_0, \{k+1\}T_0]$ ; the integration leads to

$$\begin{aligned}
 (29) \quad & P[\{x < X(t) \leq x + dx\} \cap \{kT_0 < L_1(t) \leq (k+1)T_0\}] \\
 &= \beta T_0 [Q_{k-1}(x - s) - Q_k(x - s)]h(S - x).
 \end{aligned}$$

Finally, the marginal density of  $X(t)$  is found from (29) via the total probability law by summing on  $k$ . Since

$$\sum_{k=1}^{\infty} [Q_{k-1}(x - s) - Q_k(x - s)] = 1$$

we obtain  $f_X(x) = \beta T_0 h(S - x)$ . The final result is therefore

$$(30) \quad f_X(x) = \begin{cases} [1 + H(S - s)]^{-1} h(S - x) & \text{for } s < x \leq S, \\ [1 + H(S - s)]^{-1} \left\{ q_1(S - x) + \int_0^{S-s} q_1(S - x - u) h(u) du \right\} & \text{for } x \leq s, \end{cases}$$

where  $f_X(x)$  for  $x \leq s$  has been previously obtained in (25).

### EXAMPLE III: A CONTINUOUS TIME TWO-BIN SYSTEM WITH BACKLOG AND LINEAR DEPLETION

This example is a more sophisticated version of Example I; introduction of a time lag in delivery and service from a second bin in the interim are complicating features of the present example.

Imagine that the stock on hand is stored in two bins labeled respectively A and B. Demand is met by service from bin A until the stock in bin A is exhausted. An order for stock replenishment is then placed, but delivery occurs only after a random time lag  $\tau$ . During this delay the demand is met by service from bin B, with negative stock supplies being regarded as owed to demand. When delivery is made, bin A contains a random quantity of stock,  $M$ , while bin B has a fixed (nonrandom) quantity  $s$ . Depletion from stock is at a random linear rate  $\lambda$  per unit time. It is assumed that  $M$ ,  $\tau$ , and  $\lambda$  are mutually independent within each reordering period, but not necessarily from period to period.

What is the probability density of  $X(t)$ , the total stock on hand at time  $t$ ?

#### SOLUTION

In its equilibrium state, the stock on hand is a stationary random process as depicted in Figure 3. The random times  $t_n$  of stock arrival constitute an s.p.p. and the time interval  $T$  separating two consecutive intervals is a random variable

$$(31) \quad T = M\lambda^{-1} + \tau.$$

Here  $T$  may be taken as any such interval,  $M$  the stock in bin A on delivery, and  $\lambda$  the depletion rate for this interval; this follows because the statistics of these variables are the same in each interval. We shall again require the average number of stock arrivals per unit time  $\beta$ . From (3) and (31) immediately preceding,

$$(32) \quad \beta = \frac{1}{E(T)} = [E(M)E(\lambda^{-1}) + E(\tau)]^{-1}.$$

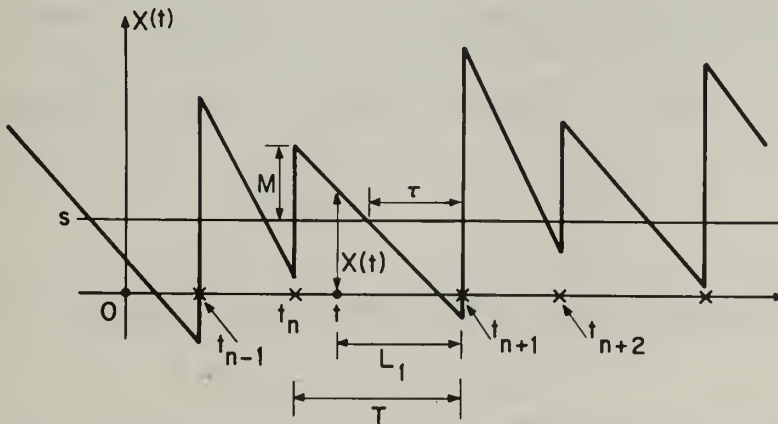


Figure 3. Two-bin inventory system with backlog and linear random depletion

The remainder of the procedure for finding  $f_X$  follows the now familiar lines of the earlier examples. Again,  $x \leq s$  and  $x > s$  are considered separately. When  $x \leq s$ , we have for  $y \leq z$

$$\begin{aligned}
 (33) \quad & P[\{x < X(t) \leq x + dx\} \cap \{y < L_1(t) \leq y + dy\} \cap \{z < \tau \leq z + dz\}] \\
 &= P[\{x < s - \lambda(z - y) \leq x + dx\} \cap \{z < \tau \leq z + dz\} \cap A_0 \cap A_1] \\
 &= P(A_1) P[\{x < s - \lambda(z - y) \leq x + dx\} \cap \{T > y\} \cap \{z < \tau \leq z + dz\}] \\
 &= \beta dy P[\{x < s - \lambda(z - y) \leq x + dx\} \cap \{z < \tau \leq z + dz\}] \\
 &= \beta dy f_\lambda \left( \frac{s - x}{z - y} \right) \frac{dx}{z - y} f_\tau(z) dz,
 \end{aligned}$$

where  $f_\lambda$  and  $f_\tau$  are the probability density functions of  $\lambda$  and  $\tau$ , respectively. For  $y > z$  the probability (33) is zero. The remainder of the computation is also like that of the earlier examples. We determine  $f_X$  as a marginal probability by integrating on  $y$  from zero to  $z$  and then on  $z$  from zero to infinity, viz.

$$\begin{aligned}
 (34) \quad f_X(x) &= \beta \int_0^\infty f_\tau(z) dz \int_0^z \frac{1}{z - y} f_\lambda \left( \frac{s - x}{z - y} \right) dy \\
 &= \beta \int_0^\infty f_\tau(z) dz \int_0^z u^{-1} f_\lambda \left( \frac{s - x}{u} \right) du.
 \end{aligned}$$

The case  $x > s$  is treated similarly. The probability in question is zero for  $y \leq z$ , and for  $y > z$

$$\begin{aligned}
 (35) \quad & P[\{x < X(t) \leq x + dx\} \cap \{y < L_1(t) \leq y + dy\} \cap \{z < \tau \leq z + dz\}] \\
 &= \beta dy P[\{x < s + \lambda(y - z) \leq x + dx\} \cap \{T > y\} \cap \{z < \tau \leq z + dz\}] \\
 &= \beta dy P \left[ \left\{ \frac{x - s}{y - z} < \lambda \leq \frac{x - s}{y - z} + \frac{dx}{y - z} \right\} \cap \left\{ \frac{M}{\lambda} + \tau > y \right\} \cap \{z < \tau \leq z + dz\} \right] \\
 &= \beta dy P \left[ \left\{ \frac{x - s}{y - z} < \lambda \leq \frac{x - s}{y - z} + \frac{dx}{y - z} \right\} \cap \{M > x - s\} \cap \{z < \tau \leq z + dz\} \right] \\
 &= \beta dy f_\lambda \left( \frac{x - s}{y - z} \right) \frac{dx}{y - z} [1 - F_M(x - s)] f_\tau(z) dz
 \end{aligned}$$

in which  $F_M$  stands for the distribution function of  $M$ . Thus the probability density for  $X(t)$  is

$$\begin{aligned}
 (36) \quad f_X(x) &= \beta[1 - F_M(x - s)] \int_0^\infty f_\tau(z) dz \int_z^\infty \frac{1}{y - z} f_\lambda\left(\frac{x - s}{y - z}\right) dy \\
 &= \beta[1 - F_M(x - s)] \int_0^\infty f_\tau(z) dz \int_0^\infty u^{-1} f_\lambda(u) du \\
 &= \beta E(\lambda^{-1}) [1 - F_M(x - s)]
 \end{aligned}$$

for  $x > s$ . The latter may be combined with  $f_X$  for  $x \leq s$  from (34) and the expression (32) for  $\beta$  to yield

$$(37) \quad f_X(x) = \begin{cases} [E(M)E(\lambda^{-1}) + E(\tau)]^{-1} \int_0^\infty f_\tau(z) dz \int_0^\infty u^{-1} f_\lambda\left(\frac{s - x}{u}\right) du & \text{for } x \leq s \\ [E(M)E(\lambda^{-1}) + E(\tau)]^{-1} E(\lambda^{-1}) [1 - F_M(x - s)] & \text{for } x > s. \end{cases}$$

In particular, if the depletion rate  $\lambda$  is a constant,  $f_\lambda(v) = \delta(\lambda - v)$  (i.e., a delta function) and consequently

$$(38) \quad f_X(x) = \begin{cases} [E(M) + \lambda E(\tau)]^{-1} \left[1 - F_\tau\left(\frac{s - x}{\lambda}\right)\right] & \text{for } x \leq s \\ [E(M) + \lambda E(\tau)]^{-1} [1 - F_M(x - s)] & \text{for } x > s. \end{cases}$$

#### EXAMPLE IV: A CONTINUOUS TIME TWO-BIN SYSTEM WITH RANDOM PARAMETER POISSON DEPLETION

The subject inventory system is a two-bin model resembling that of Example III. In place of linear depletion, the demand on the stock on hand  $X(t)$  takes the form of unit decrements which constitute a simple Poisson process; during each reordering cycle, the Poisson parameter  $\lambda$  is a random variable having the (same) probability density  $f_\lambda$ . Since we suppose that  $s$  is an integer, and  $M$  is an integer-valued random variable, the stock on hand  $X(t)$  is likewise integer-valued, as shown in Figure 4. We then ask: what is  $P[X(t) = k]$ ?

#### SOLUTION

If  $T$  is any period between successive deliveries,  $T_i$  the time interval between the  $(i-1)$ st and  $i$ 'th unit demand, and  $\tau$  is the (random) lag between ordering and delivery

$$(39) \quad T = T_1 + T_2 + \dots + T_M + \tau.$$

Because the depletion is a simple Poisson process the  $T_i$  are mutually independent, and  $f_{T_i}(w | \lambda) = \lambda e^{-\lambda w}$  for each. Hence



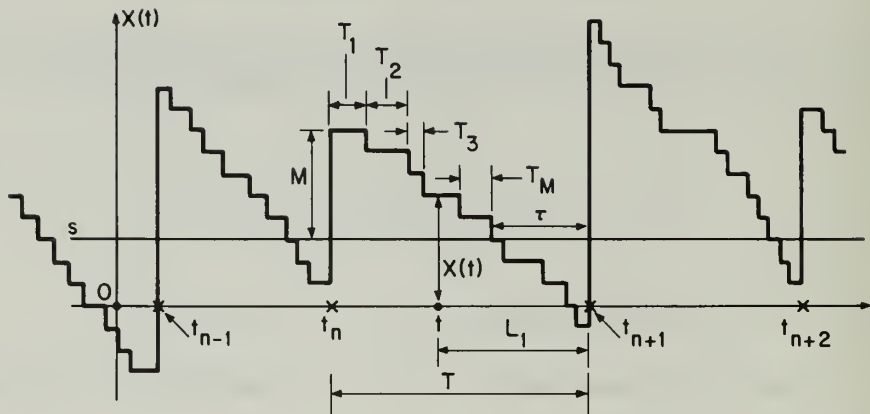


Figure 4. Continuous time two-bin inventory system with random parameter Poisson depletion

$$(40) \quad f_{T_i}(w) = \int_0^{\infty} \nu e^{-\lambda \nu} f_{\lambda}(\nu) d\nu$$

and

$$(41) \quad E(T_i) = E(\lambda^{-1}).$$

It then follows from (41) and (39) that the average rate of stock arrivals is

$$(42) \quad \beta = \frac{1}{E(T)} = [E(M)E(\lambda^{-1}) + E(\tau)]^{-1},$$

just as in Example III [cf. (32)].

To calculate  $P[X(t) = s + k]$  it is necessary to distinguish between the two cases  $k$  negative and positive. For the former, consider  $P[X(t) = s - m]$ ,  $m = 0, 1, 2, \dots$ , with  $y \leq z$ :

$$(43) \quad \begin{aligned} &P[\{X(t) = s - m\} \cap \{y < L_1(t) \leq y + dy\} \cap \{z < \tau \leq z + dz\} \cap \{w < \lambda \leq w + dw\}] \\ &= \beta dy P[E_m \cap \{T > y\} \cap \{z < \tau \leq z + dz\} \cap \{w < \lambda \leq w + dw\}] \\ &= \beta dy P[E_m \cap \{z < \tau \leq z + dz\} \cap \{w < \lambda \leq w + dw\}] \\ &= \beta dy \frac{[w(z - y)]^m \exp[-w(z - y)]}{m!} f_{\tau}(z) dz f_{\lambda}(w) dw. \end{aligned}$$

In the above,  $E_m$  denotes the event "m demands in the interval  $(t, t + z - y]$ ," and the probability in the last step is a consequence of the independence of  $\tau$  from the  $T_i$ , together with  $P[E_m \cap \{w < \lambda \leq w + dw\}] = P[E_m | \lambda = w] f_{\lambda}(w) dw$ . The first two steps in the computation (43) follow similar reasoning used in the preceding examples. For  $y > z$  the probability on the left of (43) is zero, and so the marginal probability  $P[X(t) = s - m]$  is obtained from (43) by integrating the right side on  $y$  from zero to  $z$ , and on  $z$  and  $w$  from zero to infinity. With the change of variable  $u = z - y$ , the result of these computations is

$$\begin{aligned}
 (44) \quad P[X(t) = s - m] &= \beta \int_0^\infty f_\lambda(w) dw \int_0^\infty f_\tau(z) dz \int_0^z \frac{[wu]^m}{m!} e^{-wu} du \\
 &= \beta \int_0^\infty w^{-1} f_\lambda(w) dw \int_0^\infty f_\tau(z) \left\{ 1 - \left[ \sum_{n=0}^m \frac{(wz)^n}{n!} e^{-wz} \right] \right\} dz.
 \end{aligned}$$

When  $X(t) = s + m$ ,  $m = 1, 2, \dots$ , the stock at the beginning of the reordering period must be at least  $m$ , i.e.,  $M = m + n$  with  $n$  a non-negative integer. Thus for each  $n = 0, 1, \dots$ , we consider separately

$$\begin{aligned}
 (45) \quad &P[\{X(t) = s + m\} \cap \{y < L_1(t) \leq y + dy\} \cap \{z < \tau \leq z + dz\} \cap \{w < \lambda \leq w + dw\} \cap \{M = m + n\}] \\
 &= \beta dy P[E_m \cap \{z < \tau \leq z + dz\} \cap \{w < \lambda \leq w + dw\} \cap \{M = m + n\}] \\
 &= \beta dy P[\{T_{n+1} + \dots + T_{m+n-1} < y - z \leq T_{n+1} + \dots + T_{m+n}\} \cap \{T_1 + \dots + T_{m+n} \\
 &\quad + \tau > y\} \cap \{z < \tau \leq z + dz\} \cap \{w < \lambda \leq w + dw\} \cap \{M = m + n\}] \\
 &= \beta dy P[\{T_{n+1} + \dots + T_{m+n-1} < y - z \leq T_{n+1} + \dots + T_{m+n}\} \cap \{z < \tau \leq z + dz\} \cap \\
 &\quad \cap \{w < \lambda \leq w + dw\} \cap \{M = m + n\}] \\
 &= \beta dy \frac{[w(y - z)]^m \exp[-w(y - z)]}{m!} f_\tau(z) dz f_\lambda(w) dw P[M = m + n].
 \end{aligned}$$

The above holds only for  $y > z$ ; for  $y \leq z$  the probability (45) is zero. We once more integrate on the variables  $y, z$ , and  $w$  to find the marginal probability. The integration is facilitated if we take  $u = y - z$ , thus obtaining an integral on  $u$  instead of  $y$  with new limits from zero to infinity (rather than  $z$  to infinity). The result is then

$$\begin{aligned}
 (46) \quad &P[\{X(t) = s + m\} \cap \{M = m + n\}] \\
 &= \left\{ \beta \int_0^\infty f_\lambda(w) dw \int_0^\infty f_\tau(z) dz \int_0^\infty \frac{[wu]^m}{m!} e^{-wu} du \right\} P[M = m + n].
 \end{aligned}$$

The integration on  $u$  is easily performed. Further, the total probability law can be applied to (46) by summing both sides for  $n = 0, 1, 2, \dots$ . On the right side this summation yields

$$(47) \quad \sum_{n=0}^{\infty} P[M = m + n] = P[M \geq m] = 1 - P[M \leq m - 1] = 1 - F_M(m - 1)$$

where  $F_M$  is the probability distribution function for  $M$ . When the indicated steps are carried out, we have

$$\begin{aligned}
 (48) \quad P[X(t) = s + m] &= \left\{ \beta \int_0^\infty f_\lambda(w) dw \int_0^\infty f_\tau(z) dz w^{-1} \right\} [1 - F_M(m - 1)] \\
 &= \left\{ \beta \int_0^\infty w^{-1} f_\lambda(w) dw \right\} [1 - F_M(m - 1)] = \beta E(\lambda^{-1}) [1 - F_M(m - 1)].
 \end{aligned}$$

Our calculations may then be summarized by combining the expression (48) for positive integers with (44) for negative integers, viz.,

$$(49) \quad P[X(t) = s + k] = \begin{cases} \left[ E(M) E(\lambda^{-1}) + E(\tau) \right]^{-1} \int_0^\infty f_\lambda(w) dw \\ \times \int_0^\infty f_\tau(z) \left\{ 1 - \left[ \sum_{n=0}^{-k} \frac{(wz)^{-n}}{(-n)!} e^{-wz} \right] \right\} dz & \text{for } k = 0, -1, -2, \dots \\ \left[ E(M) E(\lambda^{-1}) + E(\tau) \right]^{-1} E(\lambda^{-1}) [1 - F_M(k - 1)] & \text{for } k = 1, 2, \dots \end{cases}$$

A constant (nonrandom) depletion rate  $\lambda$  simplifies (49) to

$$(50) \quad P[X(t) = s + k] = \begin{cases} \left[ \lambda^{-1} E(M) + E(\tau) \right]^{-1} \int_0^\infty f_\tau(z) \\ \times \left\{ 1 - \left[ \sum_{n=0}^{-k} \frac{(\lambda z)^{-n}}{(-n)!} e^{-\lambda z} \right] \right\} dz & \text{for } k = 0, -1, -2, \dots \\ \left[ E(M) + \lambda E(\tau) \right]^{-1} [1 - F_M(k - 1)] & \text{for } k = 1, 2, \dots \end{cases}$$

## REFERENCES

- [1] Arrow, K. J., Karlin, S., and Scarf, H., Studies in the Mathematical Theory of Inventory and Production (Stanford Univ. Press, Stanford, Calif., 1958).
- [2] Beutler, F. J. and Leneman, O. A. Z., "The Theory of Stationary Point Processes," in Acta Math., **116**, 159-197 (1966).

- [3] Beutler, F. J. and Leneman, O. A. Z., "Random Sampling of Random Processes," *Information and Control*, 9, 325-346 (1966).
- [4] Bulinskaya, E., "Steady-State Solutions in Problems of Optimum Inventory Control," in *Theory of Prob. and Its Appl.*, 9, 502-507 (1964).
- [5] Scarf, H., Gilford, D., and Shelly, M., Multistage Inventory Models and Techniques (Stanford Univ. Press, Stanford, Calif., 1963).

\*       \*       \*





# THE STATUS AND IMPACT OF RELIABILITY METHODOLOGY\*

Gerald J. Lieberman†

Stanford University

## 1. SYSTEM CONFIGURATION AND RELIABILITY

There exist many definitions of reliability, depending upon the point of view of the user. However, they all have a common core which contains the following: Reliability,  $R(t)$ , is the probability that a device performs adequately over the interval  $[0, t]$ . In general, it is assumed that, unless repair or replacement occurs, adequate performance at time  $t$  implies adequate performance during the interval  $[0, t]$ . Although this definition is simple, the systems to which it is applied are generally very complex. In principle, it is possible to break down the system into black boxes with each black box being in one of two states, good and bad. Mathematical models of the system can then be abstracted from the physical processes and the theory of combinatorial probability utilized to predict the reliability of the system. The black boxes may be very dependent on each other. For any reasonable system, such a probability analysis generally becomes so cumbersome that it must be considered impractical. Hence, other methods are sought which either simplify the calculations or provide bounds on the reliability of the entire complex system.

### 1.1 Series Systems

In dividing the system into black boxes, it is often possible to structure it into a series configuration, which is the simplest and most common of all configurations. For a series structure, the system fails if any component of the system fails. A typical example of a series structure is a string of Christmas tree lights circa 1940.

Suppose there exists a system composed of  $n$  components connected in series. Let  $X_i$  be a binary random variable corresponding to the  $i^{\text{th}}$  component, and let the probability that unit  $i$  is successful be denoted by  $p_i = P\{X_i = 1\}$ . Let  $R$  denote the probability that the system will operate. Since a series system requires that each component perform successfully in order for the system to operate properly, the probability of success, e.g., the reliability, is given by

$$R = P\{X_1 = 1, X_2 = 1, \dots, X_n = 1\}.$$

When the usual terms for conditional probability are employed

$$R = P\{X_1 = 1\}P\{X_2 = 1 | X_1 = 1\}P\{X_3 = 1 | X_1 = 1, X_2 = 1\} \cdots P\{X_n = 1 | X_1 = 1, \dots, X_{n-1} = 1\}.$$

\*Presented at the Research and Technology Conference on Optimization/Statistics, Simulation, Information Processing, sponsored by the Office of Naval Research and the George Washington University on September 11-13, 1968.

†This work was supported in part by the Army, Navy, and Air Force under contract Nonr-225(53)(NR-042-002) with the Office of Naval Research.

In general, such conditional probabilities require careful analysis. For example,  $P\{X_2 = 1 | X_1 = 1\}$  is the probability that component 2 will perform successfully given that component 1 performs successfully. Consider a system where the heat from component 1 affects the temperature of component 2, and thereby its probability of success. The performance of these components are then dependent and the evaluation of the conditional probability is extremely difficult. If, on the other hand, the performance characteristics of these components do not interact, e.g., temperature of the component does not affect the performance of the other component, then the components can be said to be independent. The expression for the reliability then simplifies and becomes

$$R = P\{X_1 = 1\} P\{X_2 = 1\} \cdots P\{X_n = 1\} = p_1 p_2 \cdots p_n .$$

Such probabilities are often readily obtainable.

## 1.2 Parallel Systems

A parallel system is defined to be a system consisting of  $n$  components such that failure occurs if all components fail. Alternatively, the system operates if at least one of the  $n$  components performs satisfactorily. This property of parallel systems is often called "redundancy" (i.e., there are alternative components, existing within the system, to help the system operate successfully in case of a failure of one or more components). If the probability of failure of the  $i^{\text{th}}$  component is denoted by  $P\{X_i = 0\} = 1 - p_i$ , then

$$\begin{aligned} R &= 1 - P\{X_1 = 0, X_2 = 0, \dots, X_n = 0\} \\ &= 1 - P\{X_1 = 0\} P\{X_2 = 0 | X_1 = 0\} \\ &\quad \times P\{X_3 = 0 | X_1 = 0, X_2 = 0\} \cdots P\{X_n = 0 | X_1 = 0, X_2 = 0, \dots, X_{n-1} = 0\} . \end{aligned}$$

If the unit failures are independent, then this simplifies into

$$R = 1 - (1 - p_1)(1 - p_2) \cdots (1 - p_n) .$$

As an example of how a parallel structure may improve the reliability of a system consider the following example given by Shooman [10]. Shooman describes the reliability of a gyro system. He indicates that a high type military system would use three single-degree-of-freedom gyros, with each gyro being used to establish one of the three inertial reference axes,  $x, y, z$ . He assumes that the gyros cost \$15,000 each, so that the total cost is \$45,000. This can be viewed as a series system since the gyro system will work only if each of the three reference axes operate properly. Furthermore, the performance of the three gyros can be assumed to be independent, and each is known to have reliability 0.998. Hence the system reliability is given by

$$R = P\{X_1 = 1\} P\{X_2 = 1\} P\{X_3 = 1\} = (0.998)^3 = 0.994 .$$

Shooman then points out that \$45,000 is much too expensive for a commercial system. Instead, commercial systems would use three cheaper two-degree-of-freedom gyros at a cost of \$2,000

per unit, yielding a total cost of \$6,000. Furthermore, since each gyro establishes two reference directions, there are then two measurements for each reference axis. Again, the performance of each gyro reference axis is assumed to be independent, but with reliability 0.990. Thus, the system reliability can be obtained as follows. The probability that an axis, say  $z$ , operates satisfactorily is given by  $1 - (0.01)^2 = 0.9999$ . Hence, the probability that all three axes operate is

$$R = (0.9999)^3 = 0.9997,$$

which is superior to the reliability for the military system. Thus, the redundant cheaper system appears to yield a more reliable system, which is the best of all possible worlds. However, Shooman states that "It is important to note that some of the cheaper gyros will be offset by more frequent replacement and higher maintenance costs when the system is operated over a long period of time." This last statement appears to vitiate the conclusions of the entire example, but it still illustrates how redundancy can improve reliability.

### 1.3 Standby (Parallel) Model

In the parallel system just described all the components are turned on at the beginning of operation, and all units continue to perform until they fail. When a unit fails it remains in the system which continues to operate until all units fail, or the system completes its mission. In a standby system the components are connected in parallel, but do not operate at the same time. The operation of a standby system with  $(n - 1)$  spares can be described as follows: Assume that there exists a switch which can place any unit into operation in the system. At the start of operation this switch connects to component 1 and turns this unit on. Meanwhile components 2 through  $n$  are left in reserve (standby) in a turned-off condition. It is further assumed that the switch can detect when a component fails. When it senses improper operation, it switches to component 2, and then turns the unit on without interrupting the system performance. This continues through the  $n^{\text{th}}$  unit. The  $n$  channel standby system can fail only if all channels have failed after being activated. Hence,

$$\begin{aligned} R &= 1 - P\{X_1 = 0, X_2 = 0, \dots, X_n = 0\} \\ &= 1 - P\{X_1 = 0\}P\{X_2 = 0 | X_1 = 0\} \\ &\quad \times P\{X_3 = 0 | X_1 = 0, X_2 = 0\} \cdots P\{X_n = 0 | X_1 = 0, X_2 = 0, \dots, X_{n-1} = 0\}. \end{aligned}$$

Although this expression appears to look like that for the regular parallel case, its interpretation is quite different since only one component is operating at a time. Thus, unit  $j$  begins operation only after units  $1, 2, \dots, j - 1$  have failed. Hence, even when the performance characteristics of each unit do not interact so that they can be considered to be independent, an alternative expression for reliability other than that previously mentioned, is required for reliability evaluation purposes. For such independent units, an expression for the reliability at time  $t$  for a standby operation is given by

$$R(t) = P\{T_1 + T_2 + \cdots + T_n > t\},$$

where  $T_i$  is the time to failure of unit  $i$ .

#### 1.4 r Out of n Configuration

Some systems have a configuration such that the system operates if  $r$  out of  $n$  components function properly. A simple example is the cables that support a suspension bridge, i.e., the bridge will remain erect if  $r$  cables will support the bridge. If each component is identical, the reliability is given by the binomial expression

$$R = \sum_{k=r}^n \binom{n}{k} p^k (1-p)^{n-k},$$

where  $p$  is the probability of success of an individual component. If the components differ, the expression for the reliability is more complex, but can be obtained by brute force (direct enumeration).

#### 1.5 Bounds for Reliability of Series and Parallel Systems

It is evident that for systems which are in series or parallel, the expressions for the system reliability is simplified for systems which are composed of independent components. For systems which are dependent, evaluation of the reliability depends heavily on knowing conditional probabilities which are very complex. Hence, obtaining bounds for dependent systems based upon knowing the probabilities of individual components is worthwhile. In order to do this it is helpful to define the concept of associated binary random variables. This definition was introduced by Esary, Proschan, and Walkup [5], and is as follows:  $T_1, T_2, \dots, T_n$  are associated binary random variables if, and only if

$$\text{Cov}[\Gamma(T), \Delta(T)] \geq 0$$

for all pairs  $\Gamma, \Delta$  of binary, nondecreasing functions. This is the generalization of the usual definition for two random variables being associated. For example, if  $X$  and  $Y$  are two binary random variables  $\text{Cov}[X, Y] \geq 0$  implies that  $X, Y$  are associated. The bounds can now be given as follows:

If  $X_1, X_2, \dots, X_n$  are associated binary random variables then

$$R(t) = P\{X_1 = 1, X_2 = 1, \dots, X_n = 1\} \geq \prod_{i=1}^n P\{X_i = 1\} = \prod_{i=1}^n p_i.$$

This states that for a binary series system in which  $X_1, \dots, X_n$  are associated, the system reliability is at least as large as the product of the component reliabilities.

Another bound follows from the result which states that if  $X_1, X_2, \dots, X_n$  are associated binary random variables, then

$$1 - R(t) = P\{X_1 = 0, X_2 = 0, \dots, X_n = 0\} \geq \prod_{i=1}^n P\{X_i = 0\}.$$

This states that for a binary parallel system in which  $X_1, X_2, \dots, X_n$  are associated, the system reliability is no greater than one minus the product of the component failure probabilities.



Bounds for complex structures can also be obtained for associated random variables by approximating these structures with a structure having only series and parallel arrangements. Such results are also given in [5].

The results of this subsection depend heavily on the condition that the sequence of random variables is associated. Unfortunately, there is no working method for ascertaining whether or not such a system has associated components. For two components association is equivalent to having a nonnegative covariance; however, for three or more components there is no natural extension that enables one to make rational assumptions about this property for a system.

### 1.6 Practical Considerations Affecting Measurement of System Reliability

The previous sections considered reliability from the design point of view, and indicated that, for very complex systems, the measurement of system reliability as a function of component reliabilities is very cumbersome and often impractical. Bounds on the system reliability, however, can be obtained based upon a knowledge of the component reliabilities. Even such a scheme requires information about component reliabilities, and the question arises as to how to obtain such data. Unfortunately, these data are often taken from such documents as MIL-HDBK 217A [11] and the "Failure Rate Data Handbook" (FARADA). These data tend to report failure rates, and are usually compiled under a wide variety of conditions. To compensate for the latter, efforts are made to bolster the theory by incorporating stress factors and environmental factors. Unfortunately, the method of choosing these factors, and the subsequent use of results based upon these factors, is tantamount to using witchcraft. Finally, as will be noted subsequently, reporting of the failure rate, or alternatively the mean-time-to-failure, may not be sufficient to describe the data for reliability purposes.

## 2. CHARACTERIZATION OF RELIABILITIES

The previous material has described system reliability in terms of component reliability for important classes of systems. Existence of probabilities of success has been assumed, but very little has been said about estimating their values. In fact, in both government and industry contracts are often called for which require reliability at a fixed level, say 0.95, be guaranteed for systems where only one of a kind will be constructed. Hence, no operational testing can be made since, in many cases, testing is destructive; therefore, reliability has to be estimated based upon previous experience with similar type components. In such situations engineers may be reluctant to state reliabilities for components in terms of probabilities, but are often willing to give estimates for mean-time-to-failure. Hence, bounds on system reliability based upon this type of component data are useful, and the following material is motivated by this point of view.

### 2.1 Monotone Failure Rates

An appealing intuitive property in reliability is the failure rate function. Consider a component (or unit) and its associated random variable,  $T$ , time to failure. Denote the failure distribution by  $F$  and the density function by  $f$ . In terms of the previous discussion, define the random variable  $X = X(t)$  which takes on the values



$$\begin{aligned} 1 & \quad \text{if } T \geq t \\ 0 & \quad \text{if } T < t. \end{aligned}$$

Then

$$R(t) = P\{X = 1\} = p = 1 - F(t) = \int_t^{\infty} f(y) dy.$$

The failure rate function  $r(t)$  is defined for those values of  $t$  for which  $F(t) < 1$  by

$$r(t) = f(t)/R(t).$$

This function has a useful probabilistic interpretation; namely,  $r(t) dt$  represents the conditional probability that an object surviving to age  $t$  will fail in the interval  $[t, t+dt]$ . This function is sometimes called the hazard rate.

In many applications there is every reason to believe that the hazard function tends to increase because of the inevitable deterioration which occurs. Such a hazard function which remains constant or increases with age is said to have an increasing failure rate (IFR).

In some applications the hazard function tends to decrease. It would be expected to decrease initially, for instance, for materials that exhibit the phenomenon of "work hardening." Certain solid state electronic devices are also believed to have a decreasing failure rate. Thus, a hazard function which remains constant or decreases with age is said to have a decreasing failure rate (DFR).

## 2.2 Bounds on Reliability

Under either IFR or DFR assumptions, it is possible to obtain sharp bounds on the reliability in terms of moments and percentiles; in particular, such bounds can be derived from statements based upon the mean time to failure.

The failure rate function possesses some interesting properties. The failure rate distribution is completely determined by the failure rate function. In particular, it is easily shown that

$$R(t) = 1 - F(t) = \exp \left[ - \int_0^t r(\xi) d\xi \right].$$

Furthermore, a distribution is IFR (DFR) if

$$-\ln R(t) = \int_0^t r(\xi) d\xi$$

is convex (concave). As an example, consider a component whose failure distribution is given by the exponential, i.e.,

$$F(t) = P\{X < t\} = 1 - e^{-t/\theta}.$$

Thus,  $R(t)$  is given by  $e^{-t/\theta}$  and the hazard function is given by

$$r(t) = \frac{\frac{1}{\theta} e^{-t/\theta}}{\theta^{-t/\theta}} = \frac{1}{\theta}.$$

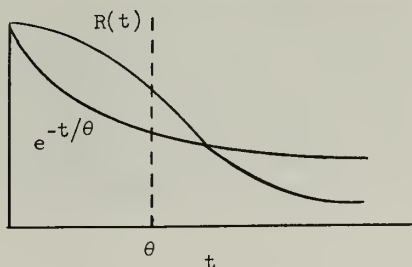
Note that the exponential has a constant failure rate and hence is both IFR and DFR. Because the exponential distribution with constant failure rate is the boundary distribution between IFR and DFR distributions, it provides natural bounds on the survival probability of IFR and DFR distributions. In particular, in [2] it is shown that if  $F$  is IFR with mean  $\theta$ ,  $R(t)$  crosses the function  $e^{-t/\theta}$  exactly once and the crossing is necessarily from above. This leads to the following bound.

If  $F$  is IFR with mean  $\theta$ , then

$$\begin{aligned} R(t) &\geq e^{-t/\theta} & t < \theta \\ &\geq 0 & t \geq \theta \end{aligned}$$

and the inequality is sharp, i.e., the exponential distribution with mean  $\theta$  attains the lower bound for  $t < \theta$ , and the degenerate distribution concentrating at  $\theta$  attains the lower bound for  $t \geq \theta$ . This can be represented graphically in Figure 2.1.

Figure 2.1. A lower bound on reliability for IFR distributions



The previous results could also provide a lower bound for  $R(t)$  whenever  $t > \theta$  provided that  $R(t)$  and  $e^{-t/\theta}$  crossed at the point  $t = \theta$ . Unfortunately this is not true and all that can be said is that the crossover is to the right of  $t = \theta$ . An upper bound [2] can be obtained from the following result. If  $F$  is IFR with mean  $\theta$ , then

$$\begin{aligned} R(t) &\leq 1 & \text{for } t \leq \theta \\ &\leq e^{-\omega t} & \text{for } t > \theta \end{aligned}$$

where  $\omega$  depends on  $t$  and satisfies  $1 - \omega\theta = e^{-\omega t}$ . Thus,  $R(t)$  can be bounded as shown in Figure 2.2.

### 2.3 Increasing Failure Rate Average

Now that bounds on the reliability of a component have been obtained, what can be said about the preservation of monotone failure rate, i.e., what structures have the IFR property

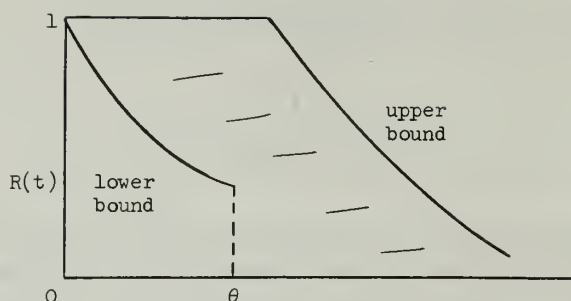


Figure 2.2. Upper and lower bounds on reliability for IFR distributions

when their individual components have this property? Series structures of independent IFR (DFR) components are also IFR (DFR). Standby systems consisting of a single IFR unit supported by  $(n - 1)$  spares is also IFR.  $k$  out of  $n$  structures consisting of  $n$  identical independent components each having an IFR failure distribution is also IFR; however, parallel structures of independent IFR components are not IFR unless they are identical components. Standby structures of DFR components are not necessarily DFR.

Thus, we see that for somewhat simple systems, there may not be a preservation of the monotone failure rate.

Instead of using the hazard rate as a means for characterizing the reliability

$$R(t) = \exp \left[ - \int_0^t r(\xi) d\xi \right],$$

a somewhat less appealing characterization can be obtained from the failure rate average function,

$$\frac{\int_0^t r(\xi) d\xi}{t} = \frac{-\log R(t)}{t}.$$

A distribution  $F$  such that  $F(0) = 0$  is called IFRA (increasing failure rate average) if and only if

$$\frac{\int_0^t r(\xi) d\xi}{t}$$

is nondecreasing in  $t \geq 0$ . A similar definition is given for DFRA. Birnbaum, Esary, and Marshall [3] showed that the IFRA class of distributions is closed under the formation of monotone binary systems. Thus, complex systems composed of components which are exponential (which is IFRA since IFR distributions are also IFRA) are also IFRA.

As with IFR systems there are bounds for systems which are IFRA. In fact, the same upper bound as given for IFR distributions is still applicable here. Barlow and Marshall [1] give the lower bound for IFRA distributions with mean  $\theta$  as follows:

$$R(t) \geq 0 \quad \text{if } t \geq \theta$$

$$R(t) \geq \min\{e^{-bt}, e^{-t/\theta}\} \quad \text{for } t < \theta,$$

where  $b$  is defined by  $e^{-bt} = b(\theta - t)$ .

## 2.4 Practical Considerations Concerning Bounds on Reliabilities

With a knowledge of the mean time to failure, bounds on the reliability can be obtained for systems having monotone failure rate functions. Similar bounds can be obtained based upon any moments of the time to failure distribution [1,2]. In any case, such an analysis only leads to bounds rather than exact reliabilities, and this will often prove to be unsatisfactory from a practical point of view. Thus, mean time to failure data is not sufficient as a characterization of the data for exact reliability calculations except when the time to failure distribution is exponential. In this case, the failure rate function is constant whereas, in general, the failure rate function depends upon time. It is evident that the reliability of components or systems depend heavily on the failure rate function, and that the assumption of constant failure rate (exponential time to failure distribution) should not be made lightly.

## 3. STATISTICAL ESTIMATION OF RELIABILITY

The previous material did not concern itself with estimating reliability except for providing bounds for components or systems. In many cases, however, data are obtained and, hence, should be utilized. Consider the simplest situation, namely, where component time to failure data are available. Let  $T$  denote the random variable, time to failure, and let  $F$  be the CDF and  $f$  the corresponding density. For this component, the expression for the reliability is given by

$$R(t) = 1 - F(t) = \int_t^{\infty} f(y) dy.$$

### 3.1 The Exponential Distribution

A very common assumption about the form of  $F(t)$  is the assumption that  $F$  is an exponential distribution, i.e.,

$$F(t) = 1 - e^{-t/\theta},$$

where  $\theta$  is easily shown to be the mean-time-to-failure ( $E(T)$ ). As was shown previously, this implies that the failure rate function, as given by  $r(t) = 1/\theta$ , is constant and independent of time. This is rather a major assumption and implies that there is no wearout; however, it is an assumption that has been made frequently, in both industry and government. Suppose that data are observed, and further, subjected to type II censoring at  $r$  out of  $n$ . Type II censoring

at  $r$  out of  $n$  refers to the situation where  $n$  items are tested simultaneously and the test is terminated when the first  $r$  out of  $n$  have failed. The time of failure of each item is recorded. This corresponds to recording the order statistics, i.e.,  $X_{(1)}, X_{(2)}, \dots, X_{(r)}$ .

Note that simple random sampling is equivalent to the case of  $r = n$ .

It is easily shown that

$$\hat{\theta} = \frac{\sum_{i=1}^r X_{(i)} + (n-r)X_{(r)}}{r}$$

is the maximum likelihood estimate for  $\theta$  so that the maximum likelihood estimate for the reliability at time  $t$  is given by

$$\hat{R}(t) = e^{-t/\hat{\theta}}.$$

Furthermore, a lower confidence bound for  $R(t)$ , with confidence coefficient  $\gamma$ , is given by

$$\exp \left[ - \frac{t \chi_{1-\gamma; 2r}^2}{2r \hat{\theta}} \right],$$

where  $\chi_{1-\gamma; 2r}^2$  is the  $1 - \gamma$  percentile of the chi square distribution with  $2r$  degrees of freedom. This is a nice result which we might try to extend to a simple series system consisting of  $n$  independent components each of which have times to failure which are exponentially distributed. Unfortunately, exact results are not easily obtainable, but a reasonable approximation due to Kramer [7] has been obtained.

Recently Sarkar [9] obtained an exact method for obtaining lower confidence bounds provided that sampling for all components is Type II censoring at  $r$  out of  $n$ .

The key to most proofs based upon the exponential distribution follows from the following relationships.

Let  $T_i$  denote the time of failure of the  $i^{\text{th}}$  item,  $X_{(i)}$  is the  $i^{\text{th}}$  order statistic,  $D_i$  is the time between  $i^{\text{th}}$  and  $(i-1)^{\text{st}}$  failure. Then

$$D_1 = X_{(1)}; D_i = X_{(i)} - X_{(i-1)} \quad \text{for } i \geq 2.$$

We shall show that  $D_i$  are independent, exponentially distributed, random variables.

Look at the joint density of all the order statistics, i.e.,

$$n! \prod_{i=1}^n f(X_{(i)}), \quad 0 \leq X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}.$$

For the exponential, this joint density becomes

$$n! \theta^{-n} e^{-\sum X_{(i)}/\theta}.$$



Now

$$X_{(1)} = D_1$$

$$X_{(2)} = D_1 + D_2$$

$$X_{(n)} = D_1 + D_2 + \dots + D_n.$$

Hence, the absolute value of the Jacobian of this transformation is one, i.e.,  $|J| = 1$ . The joint density of the  $D$ 's is then given by

$$\begin{aligned} & \frac{n!}{\theta^n} e^{-\frac{1}{\theta}[nD_1 + (n-1)D_2 + \dots + D_n]} \\ &= \left[ \frac{n}{\theta} e^{-\frac{n}{\theta}D_1} \right] \left[ \frac{n-1}{\theta} e^{-\frac{(n-1)}{\theta}D_2} \right] \dots \left[ \frac{1}{\theta} e^{-\frac{1}{\theta}D_n} \right]. \end{aligned}$$

Thus, the  $D_i$  are exponentially distributed with parameters  $\frac{\theta}{n+1-i}$ . Let  $\delta_i = (n+1-i)D_i$ .

It follows that  $\delta_i$  are independent identically distributed exponential random variables with parameter  $\theta$ . Now, the total operating time until the test terminates is equal to

$$n D_1 - (n-1) D_2 + \dots + D_n = \sum \delta_i$$

since all  $n$  live until the 1<sup>st</sup> fails,  $(n-1)$  until the second fails, and so on. Thus, the total life is

$$\begin{aligned} \sum \delta_i &= nX_{(1)} + (n-1)[X_{(2)} - X_{(1)}] + \dots + (n-r+1)[X_{(r)} - X_{(r-1)}] \\ &= \sum_{i=1}^r X_{(i)} + (n-r)X_{(r)}. \end{aligned}$$

Since  $\hat{\theta} = \sum \delta_i$  is just the sum of independent identically distributed exponential random variables it has a gamma distribution. Then,  $2r\hat{\theta}/\theta$  is chi square with  $2r$  degrees of freedom.

### 3.2 Weibull Distribution

As indicated earlier, the exponential has a serious drawback in that it has no memory or alternatively there is no wearout. Hence, alternative failure distributions should be explored.

A random variable  $X$  is said to have the Weibull distribution if the density is given by

$$f(x) = \frac{\beta}{\alpha} \left( \frac{x}{\alpha} \right)^{\beta-1} e^{-(x/\alpha)^\beta} \quad x \geq 0 \quad \alpha, \beta > 0,$$

where  $\beta$  is called the shape parameter and  $\alpha$  is called the scale parameter. Note that when  $\beta = 1$ , the Weibull reduces to the exponential.

The C.D.F. for the Weibull is given by

$$\begin{aligned} F(t) &= 1 - e^{-(t/\alpha)^\beta} & t \geq 0 \\ &= 0 & t < 0. \end{aligned}$$

Hence, the reliability is given by

$$R(t) = e^{-(t/\alpha)^\beta}.$$

The failure rate function is found from the expression

$$r(t) = \frac{f(t)}{R(t)} = \frac{\beta}{\alpha} \left( \frac{t}{\alpha} \right)^{\beta-1},$$

and is easily seen to be a monotonic function. For  $\beta > 1$ ,  $r(t)$  is a monotonically increasing function of  $t$ , and for  $\beta < 1$  is a monotonically decreasing function of  $t$ . This is shown in Figure 3.1.

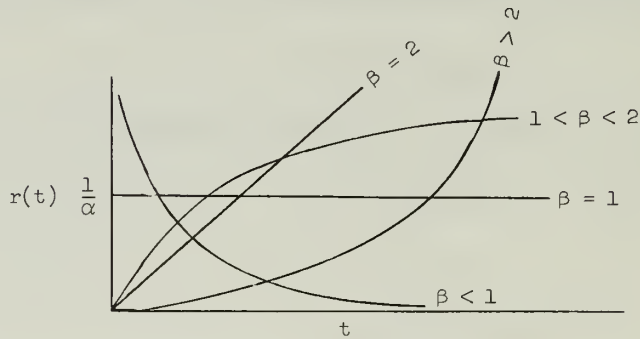


Figure 3.1. Hazard functions for the Weibull distribution

There is considerable statistical literature on estimating reliability based on the Weibull distribution. However, until a recent paper by Johns and Lieberman [6] there has been no satisfactory method presented for computing lower confidence bounds for reliability which are exact for each sample size  $n$  when both  $\alpha$  and  $\beta$  are unknown parameters. In fact, prior to the aforementioned paper, much of the literature concerned itself with finding lower confidence bounds when the shape parameter  $\beta$  was assumed to be known, a very unrealistic assumption. The Department of Defense has issued three pamphlets which are based upon this assumption of known shape parameter. It is easily shown that if one assumes that  $\beta$  is known, this problem reduces to the case of the exponential distribution. This is easily seen by operating on the transformed variable

$$Z = X^\beta,$$

which has an exponential distribution with parameter  $\theta = \alpha^\beta$ . Hence, if  $X$  has a Weibull distribution with parameters  $\alpha$  and  $\beta$  (known)

$$R(t) = e^{-(t/\alpha)^\beta} = P\{Z > t^\beta\},$$

where  $Z = X^\beta$  has an exponential distribution with parameter  $\theta = \alpha^\beta$ . All the deficiencies of the exponential assumption are really present, although somewhat obscured, when the shape parameter  $\beta$  is assumed known.

When both  $\alpha$  and  $\beta$  are assumed unknown, the traditional methods for computing lower confidence bounds lead to analytical difficulties which preclude getting exact solutions. The Johns-Lieberman paper gives the first satisfactory method for computing exact confidence bounds. This paper presents a simple method for obtaining exact lower confidence bounds for reliabilities for items whose lifetimes follow a Weibull distribution where both the scale parameter  $\alpha$  and shape parameter  $\beta$  are unknown. These confidence bounds are obtained when the data have been subjected to Type II censoring at  $r$  out of  $n$ , and are asymptotically efficient. They are exact even for small sample sizes in that they attain the desired confidence level precisely. Tables are given of exact lower confidence bounds for the reliability for sample sizes of 10, 15, 20, 30, 50, and 100 and for various values of  $r$  (including  $r = n$ ) and for various confidence coefficients  $\gamma$ . Because of the difficulty in obtaining these bounds analytically, a simulation was required for the construction of these tables. Asymptotic expressions for the lower confidence bounds which are analytically tractable are also given.

### 3.3 General Failure Rate Functions

The Weibull distribution is considered to be an important time to failure distribution since it is a "rich class," i.e., it is a two parameter family which seems to fit many sets of data. In fact, on the basis of sample data it would be virtually impossible to conclude that the time to failure distribution did not come from a Weibull distribution; however, there are other families of distributions with similar properties such as the Gamma, Log Normal, and so forth.

Another method for describing time to failure distributions is to investigate the failure rate function. It has been pointed out that

$$R(t) = 1 - F(t) = \exp \left[ - \int_0^t r(\xi) d\xi \right],$$

or alternatively,

$$F(t) = 1 - \exp \left[ - \int_0^t r(\xi) d\xi \right].$$

The simplest form of function for  $r$  is to assume that  $r$  is constant. This is the case of the exponential previously mentioned. When wear or deterioration is present, the hazard function will increase as time passes. The simplest increasing-hazard model that can be postulated is the linear function  $r(t) = Kt$  for  $t \geq 0$ . This leads to a time to failure distribution of the form

$$F(t) = 1 - e^{-Kt^2/2},$$

which is the form of the Raleigh distribution. Another type of failure rate function to consider is given by  $r(t) = Kt^m$  which is equivalent to the Weibull distribution.

In some cases the failure rate function is initially constant and then begins to increase rapidly. This could be represented by a combination of a constant and a linearly increasing hazard or an exponentially growing hazard. The hazard function

$$r(t) = Ke^{\alpha t}$$

leads to the extreme value distribution for the time to failure, i.e.,

$$F(t) = 1 - \exp[-(K/\alpha)(e^{\alpha t} - 1)].$$

Depending upon the choice of parameters, the failure rate functions usually lead to monotone functions; however, in many practical applications, the so called bath-tub function is desirable. Such a failure rate function appears in Figure 3.2. A hazard function which leads to this bath-tub shape is given by

$$r(t) = \beta/t (\alpha t^\beta)^\alpha e^{-\alpha t^\beta} \bigg/ \int_{\alpha t^\beta}^{\infty} u^{\gamma-1} e^{-u} du,$$

and the corresponding CDF is given by

$$F(t) = 1 - \frac{1}{\Gamma(\gamma)} \int_{\alpha t^\beta}^{\infty} u^{\gamma-1} e^{-u} du.$$

This is called a mixed Weibull-Gamma distribution. When  $\gamma = 1$ , the Weibull is obtained. When  $\beta = 1$ , the Gamma is obtained.

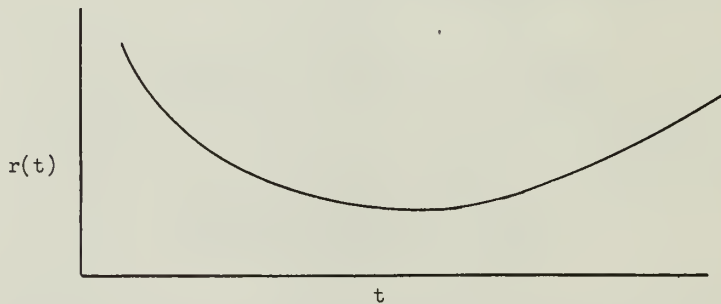


Figure 3.2. The bath-tub hazard function

### 3.4 Practical Considerations Related to the Statistical Estimation of Reliability

There exist statistical procedures for obtaining lower confidence limits on component reliability assuming a parametric form for the underlying time to failure distribution. This requires obtaining data from tests, many of which are performed in a laboratory or simulated environment. Hence, these lower confidence limits are often subject to question. In addition, the assumption of a parametric form for the underlying time-to-failure distribution requires serious study. All too often, this assumption is made without due consideration being given to its implication: assuming an exponential implies no deterioration; assuming a Weibull implies a monotone hazard function; and so on. In general, the Weibull distribution is a "rich" class of distributions having a monotone hazard function, i.e., it is a two parameter family having great flexibility. It is almost impossible to view some sample data and conclude that they do not come from the class of Weibull distributions (which includes the exponential as a special case). Of course, other two parameter families have this property.

Finally, lower confidence bounds for reliability are usually readily obtainable for components. Even for simple systems, where the form of the time-to-failure distributions are assumed known, there exist few results for finding lower confidence bounds for the overall system based upon component data.

## 4. ESTIMATION OF RELIABILITY BASED UPON ATTRIBUTES

### 4.1 Lower Confidence Limits

The previous material on statistical estimation of reliability was devoted to a discussion of time-to-failure data. If this type of data are not available, and instead binary data (success or failure) are recorded, statistical estimates of reliability can still be obtained.

Let  $X$  be a binary random variable. If the system (or component) performs satisfactorily during a test, the random variable  $X$  takes on the value 1. If a failure occurs, the random variable  $X$  takes on the value 0. Assume that

$$P\{X = 1\} = p \quad \text{and} \quad P\{X = 0\} = 1 - p$$

so that the reliability is just the probability of a successful test, i.e.,

$$R = p.$$

Thus, estimating  $R$  is equivalent to estimating the parameter,  $p$ , of a Bernoulli random variable. Suppose a sequence of  $n$  independent Bernoulli trials with parameter,  $p$ , is made and the random variables  $X_1, X_2, \dots, X_n$  are observed. It is well known that,  $\sum X_i$ , which corresponds to the total number of successes, has a binomial distribution with parameters  $n$  and  $p$ . The value of  $p$  which satisfies the expression

$$\sum_{k=\sum X_i}^n \binom{n}{k} p^k (1-p)^{n-k} = 1 - \gamma$$

is a lower  $\gamma$  confidence bound for  $p$ .



Suppose a series system consisting of  $k$  independent components is of interest. Further, tests are to be made on each of the  $k$  components, and a confidence interval estimate for the system reliability is desired. This is still an open research problem. Several approximate solutions exist, and one that is fairly good is as follows. Assume that the reliability for each component,  $p_i$ , is fairly close to 1 and the number of tests on each component,  $n_i$ , is fairly large and equal to a common value  $n$ . Then, from the Poisson approximation to the binomial  $1 - \frac{\lambda^*(D)}{n}$  is an approximate  $\gamma$  lower confidence bound for the system reliability where

$$\lambda^*(D) = \frac{1}{2} \chi_{1-\gamma; 2(D+1)}^2$$

and  $\chi_{1-\gamma; 2(D+1)}^2$  is the  $1 - \gamma$  percentile of the chi square distribution with  $2(D + 1)$  degrees of freedom and  $D$  represents the total number of failures observed, totaled over all  $k$  components. An alternative conservative approximation is given as follows:

Denote by  $\underline{p}_1, \underline{p}_2, \dots, \underline{p}_k$  the  $\gamma^{1/k}$  lower confidence bound for each component respectively (found using the binomial expression given earlier). Then  $\prod_{i=1}^k \underline{p}_i$  is a  $\gamma$  lower confidence bound for the reliability of the system.

In a recent paper, Borsting and Woods [12] present an approximate method for constructing system reliability using component failure data when the sample sizes for testing the component parts differ greatly. This is still an area where further research may prove fruitful.

## 4.2 Practical Considerations Related to Obtaining Lower Confidence

### Limits Using Attributes Data

In performing tests on complex systems, the systems are tested, modified, redesigned, and so forth depending upon the outcomes of the tests so that a true random sample is never obtained, thereby opening the theoretical results to question. In actuality, there is a reliability growth taking place which is often ignored in reporting lower confidence bounds. Reliability growth generally occurs on all programs, and should be included in reliability estimation.

## 5. UNDERLYING PHYSICAL MODELS OF RANDOM VARIABLES

Section 3 was concerned with estimating the reliability of a component or a system given the parametric form of the time-to-failure distribution. It was pointed out that an exponential time-to-failure distribution leads to a constant hazard function, which in turn implies that "no aging" takes place. This distribution can be obtained, however, by making other (equivalent) assumptions about the underlying physical process. The purpose of this section is to look at such assumptions for a variety of processes. The alternative to this approach is to look at sample data and, thereby, assume the parametric form of the distribution of the random variable, which is often an unsatisfactory procedure.

### 5.1 Exponential Failure Law

Suppose that a component fails because of the appearance of unpredictable random disturbances. These may be caused by the appearance of certain external forces like sudden changes in the environment, or the appearance of certain internal forces such as a malfunctioning part. Denote by  $X(t)$  the number of such disturbances occurring up to time  $t$ . Further,

assume that for any fixed  $t$ , the random variable  $X(t)$  has a Poisson distribution with parameter  $t/\theta$  (Poisson Process). Suppose that the component fails during the time period  $[0, t]$  if at least one such disturbance occurs. Under these assumptions, it is easily shown that the distribution of the time to failure is exponential.

## 5.2 Gamma Failure Law

Suppose the model of Section 5.1 is generalized as follows: As before, assume that the disturbances occur according to a Poisson process; now assume that the component fails whenever  $k$  or more disturbances occur during the time interval  $[0, t]$ ; thus, the CDF for the time-to-failure distribution is given by

$$F(t) = 1 - P\{T > t\}.$$

Since  $T > t$  if and only if  $(k-1)$  or fewer disturbances has occurred, it follows that

$$F(t) = 1 - \sum_{j=0}^{k-1} \frac{(t/\theta)^j e^{-t/\theta}}{j!},$$

which is the CDF of the Gamma distribution.

## 5.3 Weibull Failure Law

Assume that a component (or system) can be characterized by a chain made up of identical links in the sense that the distributions of breaking strengths of each link in the chain are the same. The chain breaks (component fails) when its weakest link fails. A stress is applied to the chain as a whole and it is assumed to be applied equally to each link. This, the probability distribution of the time to failure of such a component is just the probability distribution of the first order statistic, the minimum. If  $F$  represents the CDF for each link and  $G$  represents the CDF of the chain consisting of  $n$  links, then

$$G(t) = 1 - [1 - F(t)]^n.$$

It is easily seen that if  $F$  is Weibull then so is  $G$ . In fact, if  $F(t) \sim ct^\alpha$  for  $t$  close to zero ( $\alpha > 0$ ) then  $G$  has a limiting Weibull distribution. This includes the exponential, Weibull, Beta, and so on.

## 5.4 Log Normal Failure Distributions

A model which leads to the Log Normal failure distribution will be described in terms of a fatigue crack. Let  $X_1 < X_2 < \dots < X_n$  be a sequence of random variables such that  $X_i$  denotes the size of the fatigue crack at stage  $i$  of its growth. Assume that the crack growth  $X_{i+1} - X_i$  is proportional to the crack size  $X_i$ , e.g.,  $X_{i+1} - X_i = \lambda_{i+1} X_i$  for  $i = 0, 1, 2, \dots, n-1$ , where  $X_0$  represents the initial size of crack,  $X_n$  the final crack size, and  $\lambda_i$  are independent positive random variables. From the above expression for the crack size, it follows that

$$X_n = (1 + \lambda_n) X_{n-1} = X_0 \prod_{i=1}^n (1 + \lambda_i).$$

Thus,  $\log X_n$  is essentially the sum of independent random variables and, by invocation of the central limit theorem, is approximately normally distributed. Hence,  $X_n$  has an approximate log normal distribution.

In a recent paper Birnbaum and Saunders [4] generalized the above model with a view towards making it more realistic. This leads to a new two-parameter family of life length distributions. Their derivation follows from considerations of renewal theory for the number of cycles needed to force a fatigue crack extension to exceed a critical value.

### 5.5 Bivariate Exponential Distribution

Marshall and Olkin [8] introduced a bivariate exponential distribution which can be derived as follows: Consider three independent types of (Poisson) events occurring in time. As indicated in Section 5.1 these may represent the appearance of unpredictable random disturbance. Let  $U$ ,  $V$ , and  $W$  represent the time between disturbances of the first, second, and third types, respectively. Assume that  $U$ ,  $V$ , and  $W$  are exponential with respective parameters  $\theta_1$ ,  $\theta_2$ , and  $\theta_{12}$ . These assumptions are essentially equivalent to those made in Section 5.1. Consider a system consisting of two components connected in series, and assume that component 1 fails if a disturbance of the first or third type occurs, and component 2 fails if a disturbance of the second or third type occurs. Let  $X$  be the time to failure of the first component and  $Y$  be the time to failure of the second component so that

$$X = \min(U, W)$$

$$Y = \min(V, W).$$

It then follows that

$$\begin{aligned} P\{X > s, Y > t\} &= P\{U > s, V > t, W > \max(s, t)\} \\ &= P\{U > s\}P\{V > t\}P\{W > \max s, t\} = e^{-s/\theta_1 - t/\theta_2 - \max(s, t)/\theta_{12}}. \end{aligned}$$

The reliability at time  $t$  is then easily found, and given by

$$R(t) = P\{X > t, Y > t\} = e^{-(1/\theta_1 + 1/\theta_2 + 1/\theta_{12})t},$$

which is called the bivariate exponential distribution. Marshall and Olkin have also considered a multivariate exponential distribution which is essentially a generalization of the above arguments.

### 5.6 Practical Considerations Pertaining to Physical Models

As indicated earlier, the arbitrary choice of a parametric form of the underlying time to failure distribution is often an unsatisfactory procedure. An alternative to this is to investigate the underlying physical process which induces the time to failure distribution. This may lead to an appropriate choice of the parametric form of this distribution. Of course, no parametric form should be assumed, no matter how reasonable the derivation from physical law appears, until it is compared with actual data.

## REFERENCES

- [1] Barlow, R. E. and A. W. Marshall, "Bounds on Interval Probabilities for Restricted Families of Distributions," 5th Berkeley Symposium (1965), pp. 1229-257.
- [2] Barlow, R. E., and F. Proschan, Mathematical Theory of Reliability (John Wiley and Sons, Inc., New York, N.Y., 1965).
- [3] Birnbaum, Z. W., J. D. Esary, and A. W. Marshall, "A Stochastic Characterization of Wear-Out For Components and Systems," Ann. Math. Stat., 37, 816-825 (1966).
- [4] Birnbaum, Z. W. and S. C. Saunders, "A New Family of the Distributions," University of Washington Laboratory of Statistical Research, Technical Report No. 52 (1968).
- [5] Esary, J. D., F. Proschan, and D. Walkup, "A Multivariate Notion of Association, with Applications," Ann. Math. Stat., 38, 1466-1474 (1967).
- [6] Johns, M. V. and G. J. Lieberman, "An Exact Asymptotically Efficient Confidence Bound for Reliability in the Case of the Weibull Distribution," Technometrics, 8, 135-175 (1966).
- [7] Kramer, H. C., "One-Sided Confidence Interval for the Quality Indices of a Complex Item," Technometrics, 5, 400-403 (1963).
- [8] Marshall, A. W. and I. Olkin, "A Multivariate Exponential Distribution," J. Am. Statist. Assoc. 62, 30-44 (1967).
- [9] Sarkar, T., Unpublished Manuscript on Reliability, Stanford University (1968).
- [10] Shooman, M. L., Probabilistic Reliability — An Engineering Approach (McGraw Hill Book Co., New York, N.Y., 1968).
- [11] U.S. Government, "Reliability Stress and Failure Rate Data," MIL-HDBK-217A, Government Printing Office, Washington, D.C. (1965).
- [12] Woods, W. M., and J. R. Borsting, "A Method for Computing Lower Confidence Limits on System Reliability Using Component Failure Data with Unequal Sizes," U.S. Naval Postgraduate School Technical Report NPS55 Wo/Bg 8061 A (1968).

\* \* \*





# A NEW DERIVATION OF THE LOGISTIC DISTRIBUTION\*

Satya D. Dubey

*Department of Industrial Engineering and Operations Research  
New York University*

## ABSTRACT

Logistic distribution is widely used in describing biological, engineering, industrial, and various other types of data. In this paper it is shown that this distribution is a special case of a compound generalized extreme value distribution which is derived by compounding a generalized extreme value distribution with a gamma distribution. The paper contains several useful results relevant to these distribution functions.

## 1. GENERALIZED EXTREME VALUE DISTRIBUTION

Let the probability density function (p.d.f.) of a generalized extreme value distribution of double exponential type be represented by

$$(1.1) \quad f(x) = \alpha \exp\left(\frac{x - \mu}{\sigma}\right) \exp\left[-\alpha \sigma \exp\left(\frac{x - \mu}{\sigma}\right)\right], \quad -\infty < x < \infty, (\alpha, \sigma > 0, -\infty < \mu < \infty).$$

If we substitute  $\sigma$  for  $\alpha^{-1}$  in the above expression (1.1), it reduces to the p.d.f. of the extreme value distribution of double exponential type (Refs. [2,5]). We may obtain the moments of the above generalized extreme value distribution by computing the moment generating function (m.g.f.)  $\phi(t)$  of the random variable (r.v.)  $z = (x - \mu)/\sigma$ , where the r.v.  $x$  obeys the above generalized extreme value distribution. Thus we find  $\phi(t)$  to be

$$(1.2) \quad \phi(t) = (\alpha \sigma)^{-t} \Gamma(1 + t),$$

where  $\Gamma(1 + t)$  is the usual gamma function. From expression (1.2) we find the mean,  $EX$ , and the variance,  $\text{Var } X$ , of the above generalized extreme value distribution as

$$(1.3) \quad EX = \mu + \sigma \{\Gamma'(1) - \ln \alpha \sigma\}$$

and

$$(1.4) \quad \text{Var } X = \sigma^2 [\Gamma''(1) - (\Gamma'(1))^2],$$

where  $\Gamma'(1)$  and  $\Gamma''(1)$  are the first and the second derivatives of the gamma function  $\Gamma(x)$  evaluated at the point  $x = 1$ . These constants are

---

\*This paper was partially supported by the Army Research Office (Durham) in affiliation with the Army Munitions Command under Contract DA31-124-AROD-338 with New York University.

$$\Gamma'(1) = -\gamma \text{ (Euler's Constant)} = -0.57721 \dots$$

and

$$\Gamma''(1) = \gamma^2 + \pi^2/6,$$

respectively. Since the cumulative distribution function (c.d.f.) of the above generalized extreme value distribution is

$$(1.5) \quad F(x) = 1 - \exp \left[ -\alpha \sigma \exp \left( \frac{x - \mu}{\sigma} \right) \right],$$

we immediately find the linear relationship

$$(1.6) \quad y = A + Bx,$$

where  $y = \ln(-\ln(1 - F(x)))$ ,  $A = \ln \alpha + \ln \sigma - (\mu/\sigma)$  and  $B = 1/\sigma$ . Thus by plotting a given set of data on the proper graph paper we can decide if it is described by this generalized extreme value distribution. From expression (1.5) we get the median  $\lambda$  of this distribution to be

$$(1.7) \quad \lambda = \mu + \sigma (\ln \ln 2 - \ln \alpha \sigma).$$

Expression (1.4) provides an explicit estimator for  $\sigma$ . Also, from expressions (1.3) and (1.7), we get the relationship

$$(1.8) \quad \sigma = \frac{EX - \lambda}{\Gamma'(1) - \ln \ln 2},$$

which provides another explicit estimator for  $\sigma$ .

The generalized extreme value distribution of this note enjoys a very useful property.

Property: Let  $X_1, X_2, \dots, X_n$  be  $n$  independent, identically distributed, random variables from a generalized extreme value distribution with the parameters  $\alpha, \mu$ , and  $\sigma$  and let  $Y_n = \min(X_1, X_2, \dots, X_n)$ . Then  $Y_n$  obeys a generalized extreme value distributed with the parameters  $n\alpha, \mu$ , and  $\sigma$ . Conversely, if  $Y_n$  has a generalized extreme value distribution with the parameters  $\beta, \eta$ , and  $\delta$ , then each  $X_i (i=1, 2, \dots, n)$  obeys a generalized extreme value distribution with the parameters  $n^{-1}\beta, \eta$ , and  $\delta$ . The proof and the application of this property are given in Ref. [1].

## 2. A COMPOUND GENERALIZED EXTREME VALUE DISTRIBUTION

Assume that the parameter  $\alpha$  of (1.1) has the gamma p.d.f. of the form

$$(2.1) \quad f(\alpha) = \begin{cases} \frac{q^p \alpha^{p-1} \exp(-q\alpha)}{\Gamma(p)}, & \alpha > 0, (p, q > 0), \\ 0 & \text{elsewhere.} \end{cases}$$

Then the p.d.f. of a compound generalized extreme value distribution based on expressions (1.1) and (2.1) is

$$(2.2) \quad f(x) = \frac{pq^p \exp\left(\frac{x - \mu}{\sigma}\right)}{\left[q + \sigma \exp\left(\frac{x - \mu}{\sigma}\right)\right]^{p+1}}, \quad -\infty < x < \infty.$$

If we put  $p = 1$  and  $q = \sigma = \sqrt{3} \delta / \pi$ , ( $\delta > 0$ ), in (2.2) then it reduces to the p.d.f. of the well-known logistic distribution (Ref. [3]). Its intensity or hazard function  $\lambda(x)$  is

$$(2.3) \quad \lambda(x) = \frac{f(x)}{1 - F(x)} = \frac{p \exp\left(\frac{x - \mu}{\sigma}\right)}{q + \sigma \exp\left(\frac{x - \mu}{\sigma}\right)},$$

which can be rewritten as

$$(2.4) \quad \lambda(x) = \frac{a}{1 + b \exp\left[-\left(\frac{x - \mu}{\sigma}\right)\right]},$$

where

$$F(x) = 1 - q^p / \left(q + \sigma \exp\left(\frac{x - \mu}{\sigma}\right)\right)^p, \quad a = p/\sigma, \text{ and } b = q/\sigma.$$

Now it should be clear to the readers that  $\lambda(x)$  can be considered to be a generalized logistic curve (Ref. [4], pp. 658-661). The m.g.f.  $\phi(t)$  of this compound generalized extreme value distribution corresponding to the r.v.  $z = (x - \mu)/\sigma$  is

$$(2.5) \quad \phi(t) = p \left(\frac{q}{\sigma}\right)^t B(p - t, t + 1),$$

provided  $t < p$ . Here  $B(p - t, t + 1)$  is the usual Beta function. Applying formula (2.5), we find the mean  $EX$  and the variance  $\text{Var } X$  of this compound generalized extreme value distribution as

$$(2.6) \quad EX = \mu + \sigma \{\ln(q/\sigma) - (\Gamma'(p)/\Gamma(p)) + \Gamma'(1)\}$$

and

$$(2.7) \quad \text{Var } X = \sigma^2 \left[ \frac{\Gamma''(p)}{\Gamma(p)} - \left(\frac{\Gamma'(p)}{\Gamma(p)}\right)^2 + \Gamma''(1) - (\Gamma'(1))^2 \right],$$

where  $\Gamma'(p)$  and  $\Gamma''(p)$  are the first and the second derivatives of the gamma function  $\Gamma(x)$ , evaluated at the point  $x = p$ . For higher order moments of this distribution, we suggest the use of the formula

$$EX^r = E_{\alpha} E_X(X^r/\alpha),$$

since it is easier to compute the higher order moments of the generalized extreme value distribution, discussed in Section 1, from its moment generating function, given by expression (1.2). The actual computation is left as an exercise for interested readers.

The compound generalized extreme value distribution of this note enjoys a very useful property.

Property: Let  $X_1, X_2, \dots, X_n$  be  $n$  independent, identically distributed, random variables from a compound generalized extreme value distribution with the parameters  $p, q, \mu$ , and  $\sigma$ , and let  $Y_n = \min(X_1, X_2, \dots, X_n)$ . Then  $Y_n$  obeys a compound generalized extreme value distribution with the parameters  $np, q, \mu$ , and  $\sigma$ . Conversely, if  $Y_n$  has a compound generalized extreme value distribution with the parameters  $\beta, \eta, \delta$ , and  $\nu$ , then each  $X_i (i=1, 2, \dots, n)$  obeys a compound generalized extreme value distribution with the parameters  $\beta/n, \eta, \delta$ , and  $\nu$ . The proof and the application of this property are given in Ref. [1].

#### REFERENCES

- [1] Dubey, S. D., "Characterization Theorems for Several Distributions and Their Applications," J. Industrial Math. Soc., 16, 1-22 (1966).
- [2] Gumbel, E. J., Statistics of Extremes (Columbia University Press, New York, N.Y., 1958).
- [3] Gupta, S. S. and Shah, B. K., "Exact Moments and Percentage Points of the Order Statistics and the Distribution of the Range from the Logistics Distribution," Ann. Math. Statist., 36, 907-920 (1965).
- [4] Hald, A., Statistical Theory with Engineering Applications (John Wiley & Sons, Inc., New York, N.Y., 1952).
- [5] Parzen, E., Stochastic Processes (Holden-Day, Inc., San Francisco, Calif., 1962).

\* \* \*

ON THE THEORY OF SEMI-INFINITE PROGRAMMING AND A  
GENERALIZATION OF THE KUHN-TUCKER SADDLE POINT  
THEOREM FOR ARBITRARY CONVEX FUNCTIONS\*†

A. Charnes‡

*Northwestern University*

W. W. Cooper

*Carnegie-Mellon University*

and

K. O. Kortanek

*Cornell University*

ABSTRACT

We first present a survey on the theory of semi-infinite programming as a generalization of linear programming and convex duality theory. By the pairing of a finite dimensional vector space over an arbitrarily ordered field with a generalized finite sequence space, the major theorems of linear programming are generalized. When applied to Euclidean spaces, semi-infinite programming theory yields a dual theorem associating as dual problems minimization of an arbitrary convex function over an arbitrary convex set in  $n$ -space with maximization of a linear function in non-negative variables of a generalized finite sequence space subject to a finite system of linear equations.

We then present a new generalization of the Kuhn-Tucker saddle-point equivalence theorem for arbitrary convex functions in  $n$ -space where differentiability is no longer assumed.

INTRODUCTION

The first part of this paper is a survey of the theory of semi-infinite programming as a generalization of linear programming and convex duality theory. By the pairing of a finite dimensional vector space over an arbitrarily ordered field with a generalized finite sequence space, the major theorems of linear programming are generalized. When applied to Euclidean spaces, semi-infinite programming theory yields a dual theorem associating as dual problems minimization of an arbitrary convex function over an arbitrary convex set in  $n$ -space with maximization of a linear function in non-negative variables of a generalized finite sequence

---

\*This work was presented as an invited paper for the Joint European Meeting of the Econometric Society and the Institute of Management Sciences, Warsaw, Poland, September 1966.

†Part of the research underlying this report was supported by the Office of Naval Research projects (Contract Nonr-1228(10), Project NR 047-021 at Northwestern University, Contract Nonr-760(24), Project NR 047-48 at Carnegie-Mellon University) and for the U.S. Army Research Office--Durham, Contract No. DA-31-124-ARO-D-322 at Northwestern University, and NIH, 1 P10 ES 00098-01 at Cornell University.

‡Presently Jesse H. Jones Professor of System Science, University of Texas.



space subject to a finite system of linear equations. We present the main results on canonical closures of linear inequality systems, the Inhomogeneous Haar Theorem, and the Extended Dual Theorem for canonically closed linear inequality systems. We also discuss several regularization procedures with application to removing duality gaps which arise in compact systems not necessarily canonically closed, and with application to the differentiable convex programming example of Slater which does not satisfy any constraint qualification.

In the second part of the paper we derive a generalization of the Kuhn-Tucker saddle point equivalence theorem [10] for arbitrary convex functions defined over an open convex set in  $n$ -space which have systems of supporting hyperplanes possessing the Farkas-Minkowski property. The method of proof via the Extended Dual Theorem provides direct expressions of the Lagrange multipliers occurring in the saddle value problem as positive linear combinations of the dual evaluators (members of an appropriate g.f.s.s.) associated with any Farkas-Minkowski linear inequality equivalent of the direct problem. Our theorem is closely related to Theorem 2 of Fan-Glicksberg-Hoffman [9] where it is assumed that constraint functions possess an interior point, and where the Lagrange multipliers are given implicitly, following from the separation theorem of convex sets. Their theorem is an important generalization in another direction, namely that the underlying vector space over the real field is not necessarily topologized nor of finite dimension.

#### A SURVEY OF THE THEORY OF SEMI-INFINITE PROGRAMMING

One of the basic underlying notions of the theory of semi-infinite programming is that of a generalized finite sequence space (g.f.s.s.) —  $S$  with respect to an arbitrary index set  $I$  of arbitrary cardinality over an ordered field  $F$ . By definition,  $S$  is the vector space of all vectors  $[\lambda_i : i \in I]$  over  $F$  with only finitely many non-zero entries [4, p. 211].\* Let  $V$  be a finite dimensional vector space over  $F$ , which for simplicity we may take as  $F_m$ ,  $m$ -tuples over  $F$ . Let  $P_0, \{P_i : i \in I\}$  be a collection of vectors in  $V$ . The well-known notions of "requirements space" and "solutions space" of finite linear programming [2,3] are generalized in the following way. We call the subspace  $R$  spanned by the vectors  $P_0, \{P_i : i \in I\}$  the requirements space and define a solution set  $\Lambda$  by

$$\Lambda = \left\{ \lambda \in S : \sum_{i \in I} P_i \lambda_i = P_0, \lambda_i \geq 0 \right\}$$

is convex and we call  $S$  the "solutions space" since it contains the solution set. With these generalizations the following two theorems generalize two of the fundamental theorems of linear programming [2,3,4 (pp. 211-212)].

**THEOREM 1:** [Linear Independence by Association with Extreme Points].  $\lambda \neq 0$  is an extreme point of  $\Lambda$  if and only if the non-zero coordinates of  $\lambda$  correspond to coefficients of linearly independent vectors in  $R$ .

\* $S$  may be characterized as the set of all functions defined on  $I$  with values in  $F$ , with only finitely many values being non-zero.

**THEOREM 2:** [Opposite Sign Theorem].  $\Lambda$  is generated by its extreme points if and only if for any  $\alpha \in S$ ,  $\alpha \neq 0$ ,  $\sum_{i \in I} P_i \alpha_i = 0$  implies some  $\alpha_r$  and some  $\alpha_s$  are of opposite signs.

Notice that the underlying field  $F$  may be any arbitrarily ordered field and that in Theorem 2,  $\Lambda$  need not be bounded although generated by its extreme points. Given any solutions set  $\Lambda$  in the g.f.s.s.  $S$  and a vector  $[c_i : i \in I, c_i \in F]$ , we form the following dual semi-infinite programs [1, p. 213].

(I)

$$\min u^T P_0$$

$$\text{s.t.} \quad u^T P_i \geq c_i \quad \text{all } i \in I$$

$$u \in F_m$$

(II)

$$\max \sum_{i \in I} c_i \lambda_i$$

$$\text{s.t.} \quad \sum_{i \in I} P_i \lambda_i = P_0$$

$$\lambda \in S, \lambda \geq 0.$$

When  $F$  is the real field, the linear inequality system  $u^T P_i \geq c_i$  for all  $i \in I$  is said to be canonically closed if the set of coefficients  $\{(P_i^T, c_i) : i \in I\}$  is compact in  $R_{m+1}$  and there exists a point  $u_0$  such that  $u_0^T P_i > c_i$  for all  $i \in I$ , i.e., the system has interior points. The object of primary interest will be the nature of the solutions to the linear equality system  $u^T P_i \geq c_i$ , all  $i \in I$ . The following theorems are vital in order that the canonical closure property be an invariant of the constraint set [5, p. 114].

**THEOREM 3:** Any system of linear inequalities has a canonical closure possibly of lower dimension than the original number of variables or determines a constraint set consisting of a single point.

**THEOREM 4:** Any canonical closure of a system of linear inequalities has the same solution set as the original system.

A system of linear inequalities  $u^T P_i \geq c_i$ , for all  $i \in I$ , is said to have the Farkas-Minkowski property if the following is true. If  $u^T P \geq c$ ,  $P \in R_m$ , holds whenever  $u^T P_i \geq c_i$  for all  $i \in I$ , then there exist  $\lambda_i \geq 0$  for all  $i \in I$ ,  $\lambda_0 \geq 0$ , with at most  $n + 1$  non-zero such that

$$u^T P - c = \sum_i (u^T P_i - c_i) \lambda_i + \lambda_0 \quad \text{for all } u \in R_m.$$

Systems having the Farkas-Minkowski property are called Farkas-Minkowski systems. The importance of canonical closure, is seen in the following theorem of Haar [4, p. 213].

THEOREM 5: [Inhomogeneous Haar Theorem]. If  $u^T P_i \geq c_i$  for all  $i \in I$  is a canonically closed system, then it is a Farkas-Minkowski system.

While canonical closure is sufficient to obtain the Farkas-Minkowski property for the case of consistent programs, it is that property which leads to the following extended dual theorem for problems (I) and (II).

THEOREM 6: [Extended Dual Theorem]. If the inequality system  $u^T P_i \geq c_i, i \in I$  has the Farkas-Minkowski property when consistent, then precisely one of the following occurs:

- (i) I inconsistent,  $\sup \sum_i c_i \lambda_i = \infty$ , or  $\sup \sum_i c_i \lambda_i$  finite,
- (ii) II inconsistent,  $\inf u^T P_0 = -\infty$ ,
- (iii) I and II are both inconsistent,
- (iv) both consistent and  $\inf u^T P_0 = \sup \sum_i c_i \lambda_i = \sum_i c_i \lambda_i^*$  for some  $\lambda^* \in \Lambda$ .

#### REGULARIZATION PROCEDURES

Aside from a method of proof of Theorem 6, the semi-infinite regularization procedures provide a direct way of (a) rendering inconsistent problems into consistent ones, (b) removing duality gaps which may arise in compact linear inequality systems not necessarily canonically closed, and (c) providing Farkas-Minkowski systems for certain convex programming problems whose differential constraint functions lack constraint qualifications. Problems (I) and (II) are given the following regularizations.

$(I_R)$	$(II_R)$
$\min Mt + u^T P_0$	$\max \sum_i c_i \lambda_i - U e_m^T (\nu^+ + \nu^-)$
s.t. $t + u^T P_i \geq c_i, i \in I$	s.t. $\sum_i \lambda_i + \lambda_0 = M$
$u^T I_m \geq -U e_m^T$	$\sum_i P_i \lambda_i + I_m (\nu^+ - \nu^-) = P_0$
$-u^T I_m \geq -U e_m^T$	$\lambda, \nu^+, \nu^- \geq 0,$
$t \geq 0$	

where  $M$  and  $U$  are large positive constants. Thus, with the introduction of the additional  $t$ -variable, the linear inequality system of  $(I_R)$  has interior points. Since the original system  $u^T P_i \geq c_i, i \in I$  may be assumed compact without affecting the given constraint set, and both  $(I_R)$  and  $(II_R)$  are consistent, these dual problems are now in case (iv) of Theorem 6.

In general, for compact systems not necessarily canonically closed, Duffin-Karlovitz [8] show that duality gaps can arise. They present examples where (1) problem (I) is consistent with finite infimum and (II) is inconsistent, and (2) both (I) and (II) are consistent, but the infimum of (I) does not equal the supremum of (II). From Theorems 1 and 2 above we may interpret this phenomenon as a consequence of an inappropriate representation (or coordinate system) for the constraint set; however, it is possible to work directly with original data via a partial regularization associated with  $(I_R)$  above.

**THEOREM 7:\*** Consider the compact system, not necessarily canonically closed,

$$\begin{aligned}
 & (I') \\
 & \min \quad u^T P_0 \\
 & \text{s.t.} \quad u^T P_i \geq c_i, \quad i \in I \\
 & \quad \quad u^T I_m \geq -U e_m^T \\
 & \quad \quad -u^T I_m \geq -U e_m^T.
 \end{aligned}$$

If  $u^*$  is optimal for the original problem (I) and  $U > 0$  is such that  $u^*$  is  $(I')$ -feasible, then for the dual problem  $(II')$ , the supremum exists and equals  $u^{*T} P_0$ .

As mentioned above, the regularization procedures are also applicable to differential convex programming problems where a constraint qualification is violated. This means that the natural gradient inequality system may not characterize the constraint set. We illustrate this application with a regularization developed elsewhere [7], of the original Slater example.<sup>†</sup>

• Restating the Slater example, we have:

$$\begin{aligned}
 & (I) \\
 & \min x \\
 & \text{subject to } -(1-x)^2 \geq 0 \text{ with unique optimum } x_* = 1.
 \end{aligned}$$

Introducing a differential system of supports to contain the optimum, we obtain the equivalent problem:

$$\begin{aligned}
 & (I) \\
 & \min x \\
 & \text{subject to } 2(1-\alpha)x \geq 1-\alpha^2 \text{ for } 0 \leq \alpha \leq 2.
 \end{aligned}$$

The semi-infinite dual regularizations are as follows:

\*See [5], Theorem 4, p. 119.

<sup>†</sup>See Slater [11]; [4], p. 216; and [5], p. 119.

(I<sub>R</sub>)

$$\min Mt + x$$

$$\text{subject to } t + 2(1 - \alpha) x \geq 1 - \alpha^2, \quad 0 \leq \alpha \leq 2$$

$$x \geq -U$$

$$-x \geq -U$$

(II<sub>R</sub>)

$$\max \sum_{\alpha} (1 - \alpha^2) \lambda_{\alpha} - U\lambda^+ - U\lambda^-$$

$$\text{subject to } \sum_{\alpha} \lambda_{\alpha} = M$$

$$\sum_{\alpha} 2(1 - \alpha) \lambda_{\alpha} + \lambda^+ - \lambda^- = 1$$

$$\lambda's \geq 0.$$

It has been shown [7] that the dual optimal solution is  $x_* = \frac{2M - 1}{2M}$ ,  $t_* = \frac{1}{4M^2}$ , and

$\lambda_{x_*}^* = M$ , and all other  $\lambda's = 0$ . The common objective value is  $1 - 1/4M$ . It is interesting to note that this dual solution is optimal even if  $M$  is a non-Archimedean quantity, in which case  $1 - 1/4M$  is larger than any real number less than 1, but itself is less than 1 under the extended ordering.\*

#### THE GENERAL CONVEX PROGRAMMING DUAL THEOREM†

Assume that the direct problem is

$$\min f(u) \text{ subject to } G(u) \geq 0,$$

where  $G^T(u) = (g_1(u), \dots, g_k(u), \dots)$  is a finite vector of concave functions defined over a suitable domain which defines the convex constraint set  $K$  of the  $u$ 's, and  $f(u)$  is convex. Let  $u^T P_i \geq c_i$ ,  $i \in I$  be a system of supports for  $K$  and  $z - u^T Q_{\alpha} \geq d_{\alpha}$ ,  $\alpha \in A$  be a system of supports for  $z - f(u) \geq 0$ . Then the direct problem is equivalent to:

$$\min z \text{ subject to } z - u^T Q_{\alpha} \geq d_{\alpha}, u^T P_i \geq c_i, \alpha \in A, i \in I.$$

Thus, we have Theorem 8.

\*See [6] which develops other nonstandard semi-infinite programming problems.

†See [4], p. 217.



THEOREM 8: Assuming the Farkas-Minkowski property for this system, the extended dual theorem applies to the following dual programs:

$$\begin{array}{ll}
 \text{(I)} & \text{(II)} \\
 \min z & \max \sum_{\alpha} d_{\alpha} \eta_{\alpha} + \sum_i c_i \lambda_i \\
 \text{subject to } z - u^T Q_{\alpha} \geq d_{\alpha}, \alpha \in A & \text{subject to } \sum_{\alpha} \eta_{\alpha} = 1 \\
 & - \sum_{\alpha} Q_{\alpha} \eta_{\alpha} + \sum_i P_i \lambda_i = 0 \\
 u^T P_i \geq c_i, i \in I & \eta_{\alpha}, \lambda_i \geq 0.
 \end{array}$$

The following theorem [7] is essentially equivalent to the principle of complementary slackness, which we shall need in the next section.

THEOREM 9: Assume that the above system has the Farkas-Minkowski property and that  $u_*$  solves (I), i.e., the minimum  $z_* = f(u_*)$  is attained. Then in the dual expression, (II), for  $z_*$ , the only supports which arise are those passing through the point  $(z_*, u_*)$ , i.e., the only support planes with  $\eta_{\alpha}^* \neq 0$  and  $\lambda_i^* \neq 0$  are those for which  $z_* = u_*^T Q_{\alpha} + d_{\alpha}$  and  $u_*^T P_i = c_i$ .

#### A SADDLE POINT EQUIVALENCE THEOREM FOR CONVEX FUNCTIONS WITH SUPPORT SYSTEMS HAVING FARKAS-MINKOWSKI PROPERTY

We shall consider the convex programming problem presented above as primal problem (P):

$$\begin{array}{l}
 \text{(P)} \\
 \min f(u) \\
 \text{s.t. } G(u) \geq 0, \text{ where } G(u) = (g_1(u), \dots, g_n(u)).
 \end{array}$$

We assume that  $f(u)$  is convex on an open convex set  $D$  in  $R_m$ , and that each component  $g_k(u)$  of  $G$  is concave over  $D$ ,  $1 \leq k \leq n$ . We assume further that the constraint set  $K \subset D$ . Let  $X_{g_k} = \{(z, u) \mid z \leq g_k(u), u \in D\}$ . Then  $X_{g_k}$  and its closure  $\bar{X}_{g_k}$  are both convex in  $R_{m+1}$ . Let  $\mu_i^{(k)} z \leq u^T P_i^{(k)} - c_i^{(k)}, i \in I^{(k)}$  be any system of supporting hyperplanes for  $\bar{X}_{g_k}$ , i.e.,

$$\bar{X}_{g_k} = \left\{ (z, u) \mid \mu_i^{(k)} z \leq u^T P_i^{(k)} - c_i^{(k)}, i \in I^{(k)} \right\}.$$

It follows that each  $\mu_i^{(k)} \geq 0$ ; for otherwise if some  $\mu_i^{(k)} < 0$ , then we obtain a contradiction as follows. Fix  $\bar{u} \in D$  and let  $z_n \rightarrow -\infty, z_n < g_k(\bar{u})$ , so that  $(z_n, \bar{u}) \in \bar{X}_{g_k}$ . But this implies

$$\mu_i^{(k)} z_n \leq \bar{u}^T P_i^{(k)} - c_i^{(k)},$$

as  $z_n \rightarrow -\infty$ , which is impossible. Observe that

$$\{u \mid g_k(u) \geq 0\} \cap D = \bar{X}_{g_k} \cap \{z = 0\} = \left\{ u \mid u^T P_i^{(k)} - c_i^{(k)} \geq 0, i \in I^{(k)} \right\}.$$

LEMMA 1: Let  $k$  be given,  $1 \leq k \leq n$ . Then for all  $u \in D$ ,  $u^T P_i^{(k)} - c_i^{(k)} \geq \mu_i^{(k)} g_k(u)$  for all  $i \in I^{(k)}$ .

PROOF: Given  $k$ , let  $u$  be arbitrary in  $D$ . Then

$$(g_k(u), u) \in \bar{X}_{g_k} \implies \mu_i^{(k)} g_k(u) \leq u^T P_i^{(k)} - c_i^{(k)},$$

all  $i \in I^{(k)}$ .

Q.E.D.

Consider now  $X_{(-f)} = \{z - f(u) \geq 0, z \in R_1, u \in D\}$  and let  $\mu_\alpha z - u^T Q_\alpha - d_\alpha \geq 0, \alpha \in A$ , be any system of supports for  $X_{(-f)}$ .

LEMMA 2: For any system of supports for  $X_{(-f)}$  as given above,  $\mu_\alpha f(u) \geq u^T Q_\alpha + d_\alpha$  for all  $u \in D$  and all  $\alpha \in A$ .

PROOF: From the definition of a supporting hyperplane, it certainly follows that for each  $\alpha \in A$ ,  $\mu_\alpha z - u^T Q_\alpha - d_\alpha \geq 0$  whenever  $z - f(u) > 0$  and  $(z, u) \in R_1 \times D$ . Furthermore,  $F(z, u) = z - f(u)$  is concave over  $R_1 \times D$  and therefore by Theorem 2, (i) of Fan-Glicksberg-Hoffman [9], there exists  $\theta_\alpha \geq 0$  such that

$$\mu_\alpha z - u^T Q_\alpha - d_\alpha \geq \theta_\alpha (z - f(u))$$

over  $R_1 \times D$ . Now if  $\mu_\alpha > \theta_\alpha$ , we obtain a contradiction by taking  $z$  large negative for fixed  $\bar{u} \in D$ . Similarly, if  $\mu_\alpha < \theta_\alpha$ , we obtain a contradiction by taking  $z$  large positive for fixed  $\bar{u} \in D$ . Hence  $\mu_\alpha = \theta_\alpha$  for all  $\alpha \in A$ ,

$$\mu_\alpha f(u) \geq u^T Q_\alpha + d_\alpha$$

for all  $u \in D$ .

Q.E.D.

THEOREM 10: [Generalized Saddle Point Theorem]. Assume that there exist systems of supporting hyperplanes for  $\bar{X}_{(-f)}$  and  $\bar{X}_{g_k}$ ,  $1 \leq k \leq n$  introduced above,  $\mu_\alpha z - u^T Q_\alpha \geq d_\alpha$ , for all  $\alpha \in A$  and  $u^T P_i^{(k)} \geq c_i^{(k)}$ , for all  $i \in I^{(k)}$ , which have the Farkas-Minkowski property with respect to the  $z$ -variable. Then the infimum  $z_*$  of the primal problem (P) is finite if and only if there exists  $\sigma^* \geq 0$  such that

$$\inf_{u \in D} \{f(u) - \sigma^* G(u)\} = z_* = \max_{\sigma \geq 0} \inf_{u \in D} \{f(u) - \sigma G(u)\}.$$

If the infimum is assumed at  $u^*$ , then  $\sigma_k^* = 0$  implies that the constraint  $g_k(u) \geq 0$  is redundant in the sense that it may be removed from (P) without disturbing optimality of  $u^*$ .

REMARK: If, for example, we assume that each constraint function has interior points, then an equivalent supporting system of dimension  $m$  exists which is canonically closed. Without the assumption of interiority (or in the case of differentiability, a constraint qualification) the above system  $u^T P_i^{(k)} \geq c_i^{(k)}, i \in I^{(k)}$ , may not have a canonical closure of dimension  $m$ . Under the assumption of interiority, the theorem has been proved by Fan-Glicksberg-Hoffman [9] under a strict inequality constraint set where the underlying vector space is not necessarily topologized nor finite dimensional.

PROOF OF THEOREM 10: Assume first that  $\inf f(u)$ , subject to  $G(u) \geq 0$  is finite, say  $z_*$ . Then by the Extended Dual Theorem 6, the dual semi-infinite problems (I) and (II) below fall into case (iv), with dual optimal solution  $(\eta^*, \lambda^*)$ .

(I)

(II)

 $\min z$ 

$$\max \sum_{\alpha} d_{\alpha} \eta_{\alpha} + \sum_k \sum_i c_i^{(k)} \lambda_i^{(k)}$$

$$\text{s.t. } \mu_{\alpha} z - u^T Q_{\alpha} \geq d_{\alpha}, \alpha \in A$$

$$\text{s.t. } - \sum_{\alpha} Q_{\alpha} \eta_{\alpha} + \sum_k \sum_i P_i^{(k)} \lambda_i^{(k)} = 0$$

$$u^T P_i^{(k)} \geq c_i^{(k)}, i \in I^{(k)}$$

$$\sum_{\alpha} \mu_{\alpha} \eta_{\alpha} = 1$$

$$\eta, \lambda \geq 0.$$

By Lemma 2, for each  $\alpha \in A$ ,  $\mu_{\alpha} f(u) \geq u^T Q_{\alpha} + d_{\alpha}$  for all  $u \in D$ . Therefore, we have

$$(1) \quad f(u) = \sum_{\alpha} \mu_{\alpha} \eta_{\alpha}^* f(u) \geq \sum_{\alpha} \left( u^T Q_{\alpha} + d_{\alpha} \right) \eta_{\alpha}^*$$

for all  $u \in D$  by (II) - feasibility. Furthermore, by Lemma 1, we have

$$(2) \quad \sum_k \sum_i \lambda_i^{(k)*} \left[ u^T P_i^{(k)} - c_i^{(k)} \right] \geq \sum_k \sum_i \left[ \lambda_i^{(k)*} \mu_i^{(k)} \right] g_k(u) = \sum_{k=1}^m \sigma_k^* g_k(u) = \sigma^{*T} G(u)$$

where we define, in general,  $\sigma_k = \sum_i \lambda_i^{(k)} \mu_i^{(k)}$ . Combining (1) and (2) we obtain:

$$(3) \quad \begin{aligned} f(u) - \sigma^{*T} G(u) &\geq \sum_{\alpha} \left( u^T Q_{\alpha} + d_{\alpha} \right) \eta_{\alpha}^* - \sum_k \sum_i \lambda_i^{(k)*} \left[ u^T P_i^{(k)} - c_i^{(k)} \right] \\ &= z_* + u^T \left[ \sum_{\alpha} Q_{\alpha} \eta_{\alpha}^* - \sum_k \sum_i P_i^{(k)} \lambda_i^{(k)*} \right] = z_*, \end{aligned}$$

for all  $u \in D$ . The chain of inequalities (3) uses both dual equality of functionals and dual feasibility. Therefore,

$$\inf_{u \in D} \{f(u) - \sigma^* T G(u)\} \geq z_*$$

The remainder of the proof may follow the technique of proof of Theorem 2 of Fan-Glicksberg-Hoffman [9], p. 621, which uses properties of the infimum, and which we include here for completeness. Since  $z_*$  is an infimum, given any  $\delta > 0$ , there exists  $u_\delta \in K$  such that  $z_* < f(u_\delta) < z_* + \delta$ . Given any  $\sigma^T \geq 0$ , let  $\alpha = \inf_{u \in D} \{f(u) - \sigma^T G(u)\}$ . If  $\alpha$  is finite, then

$$z_* + \delta > f(u_\delta) \geq \alpha + \sigma^T G(u_\delta) \geq \alpha$$

since  $G(u_\delta) \geq 0$ . If  $\alpha = -\infty$ , we still have  $z_* + \delta \geq \alpha$ . Since  $\delta$  is arbitrary, it follows that  $z_* \geq \alpha$ . Hence

$$\max_{\sigma \geq 0} \inf_u \{f(u) - \sigma^T G(u)\} \leq z_* \leq \inf_{u \in D} \{f(u) - \sigma^* T G(u)\}.$$

Hence

$$\max_{\sigma \geq 0} \inf_{u \in D} \{f(u) - \sigma^T G(u)\} = z_* = \inf_{u \in D} \{f(u) - \sigma^* T G(u)\}.$$

On the other hand, if the maximin condition holds over  $D$  and  $\sigma \geq 0$ , then  $f(u)$  is bounded below on  $K$ , i.e.,

$$f(u) \geq z_* + \sigma^* T G(u) \geq z_*$$

for all  $u \in K$  which implies that  $\inf f(u)$  is finite, and hence by our above argument, equal to

$$\max_{\sigma \geq 0} \inf_{u \in D} \{f(u) - \sigma^T G(u)\}.$$

In the event that the infimum is assumed at  $u^* \in K$ , the above inequalities (1), (2), and (3), together with  $G(u^*) \geq 0$  yield the following chain:

$$f(u) - \sigma^* T G(u) \geq z_* = f(u^*) \geq f(u^*) - \sigma^T G(u^*)$$

for all  $u \in D$ ,  $\sigma^T \geq 0$ . Furthermore, if some  $\sigma_k^* = 0$ , then since

$$\sigma_k^* = \sum_i \lambda_i^{(k)*} \mu_i^{(k)}$$

and  $\lambda_i^{(k)}, \mu_i^{(k)} \geq 0$ , it follows that  $\lambda_i^{(k)*} \mu_i^{(k)} = 0$  for all  $i \in I^{(k)}$ ; however, if  $\lambda_i^{(k)*} > 0$ , then by

(complementary slackness) Theorem 9, it follows that  $u^{*T} P_i^{(k)} = c_i^{(k)}$ . But in view of Lemma 1, this means that  $\mu_i^{(k)} > 0$ , for otherwise  $\mu_i^{(k)} = 0$ , and we obtain the contradiction that  $u^{*T} P_i^{(k)} - c_i^{(k)} \geq 0$  for all  $u \in$  open convex set  $D$  and  $u^{*T} P_i^{(k)} - c_i^{(k)} = 0$  at  $u^*$ , an interior point of  $D$ . Hence  $\sigma_k^* = 0 \implies \lambda_i^{(k)*} = 0$  for all  $i \in I^{(k)}$ , which means that the set of inequalities  $u^{*T} P_i^{(k)} \geq c_i^{(k)}$ ,  $i \in I^{(k)}$  may be removed without affecting (I) and (II) feasibility. Moreover, the equality of dual objective functions is still maintained which means that  $u^*$  is still optimal for the new primal problem. This completes the proof of Theorem 10.

## REFERENCES

- [1] Arrow, K. J., L. Hurwicz, and H. Uzawa, "Constraint Qualifications in Maximization Problems," *Nav. Res. Log. Quart.*, 8, 175-191 (June 1961).
- [2] Charnes, A., and W. W. Cooper, "The Strong Minkowski Farkas-Weyl Theorem for Vector Spaces over Ordered Fields," *Proc. Natl. Acad. Sci. U.S.*, 44, 914-916 (Sept. 1958).
- [3] Charnes, A. and W. W. Cooper, Management Models and Industrial Applications of Linear Programming (J. Wiley and Sons, New York, 1961), Vols. I and II.
- [4] Charnes, A., W. W. Cooper, and K. Kortanek, "Duality in Semi-Infinite Programs and Some Works of Haar and Caratheodory," *Management Science*, 9, 209-228 (Jan. 1963).
- [5] Charnes, A., W. W. Cooper, and K. Kortanek, "On Representations of Semi-Infinite Programs Which Have No Duality Gaps," *Management Science*, 12, 113-121 (Sept. 1965).
- [6] Charnes, A., W. W. Cooper, and K. O. Kortanek, "On Some Nonstandard Semi-Infinite Programming Problems," Technical Report No. 45, Department of Operations Research, Cornell University (Mar. 1968).
- [7] Charnes, A., W. W. Cooper, and K. Kortanek, "Semi-Infinite Programming, Differentiability and Geometric Programming: Part II," *Aplikace Matematicky*, 14, No. 1 (1969), Czechoslovakia.
- [8] Duffin, R. J. and L. A. Karlovitz, "An Infinite Linear Program with a Duality Gap," *Management Science*, 12, 122-134 (Sept. 1965).
- [9] Fan, K., I. Glicksberg, and A. J. Hoffman, "Systems of Inequalities Involving Convex Functions," *Proc. Amer. Math. Soc.* (June 1957).
- [10] Kuhn, H. W. and A. W. Tucker, "Non-Linear Programming," *Proc. 2nd Berkeley Symp. Math. Stat. and Prob.*, J. Neyman, editor (U. Calif. Press, Berkeley, Calif., 1951), pp. 481-492.
- [11] Slater, M., "Lagrange Multipliers Revisited: A Contribution to Non-Linear Programming," Cowles Commission Paper, Math. No. 403, New Haven (Nov. 1950).





Hermann Enzer

U.S. Air Force Institute of Technology

## ABSTRACT

In this paper an attempt is made to derive two interval utility measures for public goods. The domain of the measures is partitioned into priority classes in such a way that each class contains an element, called a kernel, which can be moved by parameter variation to the immediately preceding class. The "expansion path" from priority class to priority class is determined by an optimality criterion. The measures themselves are based on optimality conditions that are similar to the familiar first-order conditions of consumer analysis. One is a product measure and the other one is additive.

1. Ordinalists have held that in the area of consumer choice there is no need to continue the search for a cardinal (= interval) utility measure since the derivation of a consumer's demand curve merely requires the existence of an ordinal utility function; for example, see Ref. [2]. But when the commodities involved are collective goods, it is often necessary to make statements that presuppose a cardinal utility function. The purpose of this paper is to discuss two measures which may be useful when the commodities are weapon systems.

There have been some attempts in the Department of Defense to measure utility. For instance, the Air Force has used the so-called Churchman-Ackoff method, Ref. [1], which is a measure that lies between the ordinal and the interval scales. It is derived something like this: Assume, e.g., that four items have been ordered according to preference. Then numbers between zero and one are assigned roughly indicating the order of the preferences. Suppose that for  $X_1 > X_2 > X_3 > X_4$ , one chooses:

$X_1$	$X_2$	$X_3$	$X_4$
1.00	0.85	0.75	0.20

Next, one determines whether  $X_1$  is preferred to the combination  $X_2$ ,  $X_3$ , and  $X_4$ . If  $X_1$  is preferred, then the numbers are adjusted to reflect this situation, such as:

\*I wish to acknowledge the helpful criticism of Professor N. Georgescu-Roegen who was kind enough to read a draft of the manuscript. I am also indebted to the referee for suggesting better formulations of several concepts. Naturally, I alone am responsible for any remaining shortcomings.

†Any opinions which are expressed in this paper are the sole responsibility of the author and do not necessarily reflect the views of the organization to which he belongs.

$X_1$	$X_2$	$X_3$	$X_4$
1.00	0.65	0.20	0.10

Then one inquires whether  $X_2$  is preferred to the combination  $X_3, X_4$ . If it is, the adjusted numbers are retained. But if the combination is preferred to  $X_2$ , the numbers are changed again to, say:

$X_1$	$X_2$	$X_3$	$X_4$
1.00	0.25	0.20	0.10

Obviously, many other number combinations would have accomplished the same objective so that this measure is an ordering with certain constraints placed upon the distances.

It is apparent that the above procedure presupposes additivity,\* i.e., that there exists a utility function such that

$$U(X_1, X_2, \dots, X_i) = U(X_1) + U(X_2) + \dots + U(X_i).$$

The likelihood that there exists such a utility function in any concrete case is small, perhaps even "infinitely small"† when the elements of a commodity space are weapon systems.

Other measures, such as the von Neumann-Morgenstern measure [6], or the Pareto-Lange approach of comparing utility differences [5,7], do not appear to be promising alternatives either. It is quite certain that several of their postulates are not satisfied in the defense area. A closer approximation to an appropriate measure than the one attempted in this paper, however, may very well involve a von Neumann-like approach since there usually exist a number of weapon systems which are chosen under conditions of uncertainty. But in this paper it is assumed that all projects are certain projects.

We shall discuss the domain on which the utility functions are defined in the next section and derive the measures in Section 3.

2. First let us state the assumptions that are made concerning the domain of the utility functions and then briefly discuss each of them:

Assumption 1:  $X$  denotes a set of weapon systems. Each member of the set is described by a vector  $(p, t, q)$ , where  $p$  denotes a performance index,  $t$  the availability time, and  $q$  the quantity. Henceforth, we let  $X$  be the set of such vectors.

Assumption 2: There is a priority hierarchy of  $k$  levels, from level 1 to level  $k$ , and each  $(p, t, q) \in X$  can be assigned a unique priority level.  $P(p, t, q) = i$  denotes that  $(p, t, q)$  is assigned the priority level  $i$  by the priority function  $P$ .

\*Such assumptions are not new. They have been employed, e.g., by Fisher and Frisch; see P. A. Samuelson, *Foundations of Economic Analysis* (Harvard University Press, Cambridge, Mass., 1947), p. 174 ff. For a list of recent contributions in this area, see P. C. Fishburn, "A Note On Recent Developments in Additive Utility Theories for Multiple-Factor Situations," *Operations Research*, 14, 1143-1148 (1966).

†Samuelson, op. cit., p. 182.

Assumption 3:  $P(p, t, q) > P(p', t', q')$  implies that  $(p, t, q) > (p', t', q')$ , where  $>$  means "is preferred to."

Assumption 4: For each  $i > 1$  there exists a  $(p, t, q) \in X$  and  $p' < p, t' > t, q' < q$  such that —

- a.  $P(p, t, q) = i$ ;
- b.  $P(p^*, t^*, q^*) = i-1$  when  $(p^*, t^*, q^*) \in \{p, p'\} \times \{t, t'\} \times \{q, q'\}$  and  $(p^*, t^*, q^*) \neq (p, t, q)$ ;
- c. the set  $\{p, p'\} \times \{t, t'\} \times \{q, q'\}$

of eight elements is weakly ordered by  $\leq$  ("is not preferred to") and contains a unique least preferred element (as well as a unique most preferred element).

As far as Assumption 1 is concerned, the triple performance, time, and quantity describes a weapon system concisely and completely, at least from a static point of view. These parameters are familiar numbers to a decision maker, increasing the likelihood that any statements which he makes about parametric variations will be relatively reliable.

It is assumed in Assumption 2 that any set  $X$  can be partitioned into equivalence classes. It is not part of the definition that two elements that are in the same class must be "comparable" in terms of their parameters. The criteria which determine a partitioning may be a mixture of military, political, economic, technical, and social considerations, and also depend on the composition and size of the set  $X$ . The important point is that the criteria must be extensive enough to enable a decision maker to place each element in an equivalence class. It is believed that in the defense environment the phrase "priority level" is more suggestive and meaningful than the usual "preferred or indifferent to." Indeed, the notion of a priority hierarchy may have assumed in time a certain operational meaning with DOD decision makers so that an analyst may not have to concern himself with the criteria that are used in the formation of priority classes.

Assumption 3 points out that we are primarily interested in the priority classes rather than in the individual elements (with one exception to be discussed in the next paragraph). It should be noted that the converse of Assumption 3 does not necessarily hold, nor do we assert that  $P(p, t, q) = P(p', t', q')$  implies that  $(p, t, q)$  "is indifferent to"  $(p', t', q')$ . In order to reduce difficulties of perception indifference classes are not defined at all.\*

In Assumption 4 it is assumed that each priority class, except the first one, contains an element called the kernel of a priority class which, for some parameter variations, would be placed in the next lower class. Therefore, adjacent priority classes can be connected by means of hypothetical parameter variations. The "direction" or the "expansion path" from class to class is determined by assuming that there exists a least preferred change. For example, suppose that the variations  $(-\Delta p, 0, 0)$  and  $(0, \Delta t, 0)$  move an element from class  $i$  to class  $i-1$ . A decision maker who prefers  $(-\Delta p, 0, 0)$  to  $(0, \Delta t, 0)$ , would rather have less performance than

\*One approach to deal with problems of perception is to introduce explicitly a "perception interval," called a psychological threshold, and then define a function on this interval stating a probability that one element will be preferred over another. See N. Georgescu-Roegen, "The Pure Theory of Consumer's Behavior," Article 1 in Part II of [3]; R. D. Luce, "A Probabilistic Theory of Utility," *Econometrica*, 1958; D. Scott and P. Suppes, "Foundational Aspects of Measurement," *Journal of Symbolic Logic*, 1958; C. Villegas, "On Qualitative Probability  $\sigma$ -Algebras," *Annals of Mathematical Statistics*, 1964.

late delivery. In other words, the set of changes is not an indifference set. It is assumed that even when a decision maker is nearly indifferent among changes, he can indicate a tendency toward a least preferred change.

3. In order to derive the two utility measures we shall make the following assumptions:

Assumption 5:

- a. The Government attempts to allocate resources efficiently.
- b. In practice resources are allocated according to conditions (2) or (3) below.
- c. The ratios (2) or (3) are satisfied for all finite changes  $\{p' - p\} \times \{t' - t\} \times \{q' - q\}$  defined in Assumption 4.

Let us define an ordinal utility function on the elements of the set  $X$ :

$$V = V(p_1, t_1, q_1, p_2, t_2, q_2, \dots, p_n, t_n, q_n)$$

and assume a budget constraint

$$B = C_1(p_1, t_1, q_1) + C_2(p_2, t_2, q_2) + \dots + C_n(p_n, t_n, q_n),$$

where  $C_i$ , the cost of system  $i$ , depends on  $p_i$ ,  $t_i$ , and  $q_i$ . In addition, there usually exist side conditions, such as

$$F(p_1, p_2, \dots, p_n) = 0,$$

$$G(t_1, t_2, \dots, t_n) = 0,$$

and

$$H(q_1, q_2, \dots, q_n) = 0,$$

which have to be observed. Then, if one assumes that utility is maximized subject to the constraints, the necessary conditions for an optimal allocation of resources are

$$(1) \quad \frac{\partial V / \partial p_1}{\partial B / \partial p_1} = \frac{\partial V / \partial t_1}{\partial B / \partial t_1} = \frac{\partial V / \partial q_1}{\partial B / \partial q_1} = \dots = \frac{\partial V / \partial q_n}{\partial B / \partial q_n}.$$

It is clear that an optimal solution may call for fewer than  $n$  commodities. It should also be understood that the above ratios are theoretical ideals and that the actual functions have many critical points at which derivatives are not defined. But the important point is to recall that resources can be allocated optimally on the basis of an ordinal utility function. It is difficult to say how successful the Department of Defense is in attaining optimality, but one can presume that the Department is earnestly attempting to achieve efficient resource allocations.

But even if one takes an optimal resource allocation for granted, it is hard to believe that the actual menu of commodities is in a "realistic" neighborhood of the normative optimum. It seems more plausible to assume that decisions are based on relative and not on absolute



terms. Particularly in the case of collective goods where resources are treated to some extent "at arm's length," an allocation on the basis of relative terms often is the only way by which a decision can be reached. Relative values are surrogates for knowledge and experience in evaluating commodities. A cost tag of \$1 billion means little to most men and if a change in cost of \$1 million occurs, it seems natural to view this change relatively. Therefore it is suggested that the necessary conditions

$$(2) \quad \frac{\frac{\partial V/V}{\partial p_1}}{\frac{\partial B/B}{\partial p_1}} = \frac{\frac{\partial V/V}{\partial t_1}}{\frac{\partial B/B}{\partial t_1}} = \frac{\frac{\partial V/V}{\partial q_1}}{\frac{\partial B/B}{\partial q_1}} = \dots = \frac{\frac{\partial V/V}{\partial q_n}}{\frac{\partial B/B}{\partial q_n}}$$

based on the transformation  $U = \log V$  and  $C = \log B$ , or

$$(3) \quad \frac{\frac{\partial V}{\partial p_1}}{\frac{\partial B/B}{\partial p_1}} = \frac{\frac{\partial V}{\partial t_1}}{\frac{\partial B/B}{\partial t_1}} = \frac{\frac{\partial V}{\partial q_1}}{\frac{\partial B/B}{\partial q_1}} = \dots = \frac{\frac{\partial V}{\partial q_n}}{\frac{\partial B/B}{\partial q_n}}$$

based on  $U=V$  and  $C=\log B$ , are more accurate than the ratios (1). It is clear that other relative transformations could be used, but percentage expressions are intuitively very appealing. Although it is somewhat dangerous to use one's introspection and judge whether one assumption appears to be more reasonable than another, it is a procedure open to a social scientist.

Which set of conditions is more appropriate can be determined only by experimentation. Since conditions (2) are elasticities, it is possible that they are more realistic than conditions (3), but we shall derive a utility measure for each.

In Assumption 5c it is assumed that either set of ratios is satisfied for the finite changes that shift a kernel to the preceding class. The problem with such a requirement is that one does not know how finely partitioned a set  $X$  must be in order to justify a linearity assumption between adjacent priority classes. There is no such thing as a "natural" partition. One can only suggest that a partition should be so fine that a decision maker is just able to perceive the distinction between two adjacent priority classes.\* To obtain such a partition, one may have to fill some of the "holes" between adjacent classes either by parameter manipulation of elements already in  $X$ , or by adjoining some artificial elements to  $X$ .

Let us now select the  $3(k-1)$  ratios from (2) which designate the kernel elements and renumber them:

$$(4) \quad a = \frac{\frac{\partial V/V}{\partial p_2}}{\frac{\partial B/B}{\partial p_2}} = \frac{\frac{\partial V/V}{\partial t_2}}{\frac{\partial B/B}{\partial t_2}} = \frac{\frac{\partial V/V}{\partial q_2}}{\frac{\partial B/B}{\partial q_2}} = \dots = \frac{\frac{\partial V/V}{\partial q_k}}{\frac{\partial B/B}{\partial q_k}}.$$

\*The "marginal preference" of Armstrong or the "minimum sensible" of Edgeworth (W. E. Armstrong, "The Determinateness of the Utility Function," *The Economic Journal* (1939); F. Y. Edgeworth, *Mathematical Psychics* (London, 1881), p. 8).

Suppose, e.g., that the kernel of the second class becomes an element of the first class by a finite change in performance. We can write

$$a = \frac{\partial V_2/V_2}{\partial B_2/B_2} = \frac{(V_1 - V_2)/V_2}{\Delta B_2/B_2}$$

or

$$(5) \quad V_2 = V_1 \left( \frac{1}{1 + a \Delta B_2/B_2} \right),$$

where  $\Delta B_2/B_2$  denotes the relative decrease in cost associated with the decrease in performance that moves the kernel to the first class. If we choose an "origin,"  $V_1$ , and a unit of measurement,  $a$ , we may view (5) as an "interval" measure and let  $V_2$  denote the utility of the second priority class. We do not assume that the function  $P$  used in Assumption 2 is the utility function  $V$ .

Suppose that the kernel of the third class is put in the second class on the basis of a time change:

$$a = \frac{\partial V_3/V_3}{\partial B_3/B_3} = \frac{(V_2 - V_3)/V_3}{\Delta B_3/B_3}$$

so that

$$V_3 = V_2 \left( \frac{1}{1 + a \Delta B_3/B_3} \right),$$

or, in general,

$$(6) \quad V_{i+1} = V_i \left( \frac{1}{1 + a \Delta B_{i+1}/B_{i+1}} \right), \quad (i=1, 2, \dots, k-1)$$

$$(7) \quad V_k = V_1 \prod_{i=1}^{k-1} 1/(1 + a \Delta B_{i+1}/B_{i+1}).$$

Since  $a \Delta B_{i+1}/B_{i+1}$  is nonpositive, some care is necessary in selecting an appropriate  $a$  for which the ratios (6) remain positive. Obviously,  $V_{i+1} > V_i$ .

If there are several changes which move a kernel to the preceding class, the change which is used in (6) is the least preferred change defined in Assumption 4c. The utility measure is therefore defined on an optimal expansion path and the assigned utilities are consistent with respect to an optimality criterion. It may happen that a kernel is assigned to the preceding class not on the basis of a "pure" change, but through a combination of two or three changes. In this case one would use, say,

$$V_{i+1} = V_i \left( \frac{1}{1 + a\Delta B_{i+1}/B_{i+1}} \right)_1 \left( \frac{1}{1 + a\Delta B_{i+1}/B_{i+1}} \right)_2 \left( \frac{1}{1 + a\Delta B_{i+1}/B_{i+1}} \right)_3$$

where the subscripts denote the performance, time, and quantity directions, respectively.

If, instead of the ratios (2), we base a measure on the ratios (3), we have first for the  $k-1$  kernels:

$$(8) \quad b = \frac{\frac{\partial V}{\partial p_2}}{\frac{\partial B/B}{\partial p_2}} = \frac{\frac{\partial V}{\partial t_2}}{\frac{\partial B/B}{\partial t_2}} = \frac{\frac{\partial V}{\partial q_2}}{\frac{\partial B/B}{\partial q_2}} = \dots = \frac{\frac{\partial V}{\partial q_k}}{\frac{\partial B/B}{\partial q_k}}.$$

The equations corresponding to (5), (6), and (7) are:

$$(9) \quad V_2 = V_1 - b(\Delta B_2/B_2),$$

$$(10) \quad V_{i+1} = V_i - b(\Delta B_{i+1}/B_{i+1}),$$

$$(11) \quad V_k = V_1 - b \sum_{i=1}^{k-1} \Delta B_{i+1}/B_{i+1}.$$

The product  $b\Delta B_{i+1}/B_{i+1}$  is nonpositive and  $V_{i+1} > V_i$ . As before, we can choose an origin and a unit of measurement and take the changes defined by Assumption 4c.

Suppose that the ratios (8) are normalized so that  $\partial V = c$  for all  $i$ , where  $c$  is some constant. Then  $\partial B/B = d$ , say, for all  $i$ . If we add all the  $\Delta V_i$  that are actually used in the utility measure, we obtain

$$(12) \quad V_1 - V_k = (k-1)c,$$

and if we multiply the corresponding  $\Delta B_i/B_i$ , we have

$$\prod_{i=1}^{k-1} \Delta B_{i+1}/B_{i+1} = d^{k-1},$$

or

$$(13) \quad \prod_{i=1}^{k-1} (B_{i+1} + \Delta B_{i+1})/B_{i+1} = (1+d)^{k-1},$$

if we add one to both sides before multiplying.

Solving (12) and (13) for  $k-1$ , we derive

$$V_k = V_1 - b(d/\log(1+d)) \log \prod_1^{k-1} \left[ (B_{i+1} + \Delta B_{i+1})/B_{i+1} \right],$$

since  $c = bd$ . Letting  $\beta = (d/\log(1+d))$ ,

$$(14) \quad V_k = V_1 - b\beta \log \prod_1^{k-1} \left[ (B_{i+1} + \Delta B_{i+1})/B_{i+1} \right],$$

an expression which has some resemblance with Bernoulli's measure of moral expectations

$$V = C \log [(Y + \Delta Y)/Y].$$

As noted before, which of the two measures is a better description of reality can only be determined by experimentation.

It is possible that a priority class has more than one element which meets the definition of a kernel. In this case it is necessary to select a median or a modal kernel allowing us to speak of the kernel of the  $i$ -th priority class.

There is also the consistency problem, namely, the possibility that the distance between two non-adjacent priority classes may differ from the combined distances of the adjacent classes lying between the two selected classes. This difficulty is avoided by defining the distance between two non-adjacent classes as the product (sum) of the distances of the adjacent classes that lie between the two classes to be compared. Such a definition has some undesirable consequences.

First, in order to obtain a utility index for a class, numbers have to be assigned to all preceding classes. But since one is usually interested in the relative values of all elements of a set  $X$ , this does not seem to be much of a drawback. Secondly, and probably more serious, the measure may be a function of the number of classes into which  $X$  is partitioned. Suppose, e.g., that a new class is adjoined to  $X$ . Then the relative utilities assigned to any two classes may differ after the new class has been added even though the relation between the classes remains unchanged. This is rather serious on first sight, but there are two considerations which tend to weaken the importance of this point. The first is that the measures which have been presented are interval measures and not absolute measures. This means that intervals, or distances between classes, can be compared. If pairs of adjacent classes are compared only, the relative values do not change after the addition of new classes. Only in the case when the intervals to be compared are enlarged, does one encounter the difficulty mentioned above. The second point is that the alternative to the approach followed here is another assumption which forces consistency upon a decision maker. After considering several possible candidates, the device of defining a distance as the product (sum) of the parts seemed to be the lesser of two evils provided, of course, one wishes to remain as realistic as possible.

## BIBLIOGRAPHY

- [1] Ackoff, R. L., Scientific Method; Optimizing Applied Research Decisions (John Wiley and Sons, New York, N. Y., 1962), pp. 87-88.
- [2] Allen, R. G. D., "A Note on the Determinateness of the Utility Function," *The Review of Economic Studies*, II, 155-158 (1934-1935).
- [3] Georgescu-Roegen, N., Analytical Economics (Harvard University Press, Cambridge, Mass., 1966).
- [4] Georgescu-Roegen, N., "Measure, Quality, and Optimum Scale," *Sankhya*, Series A, 27, 39-64 (1965).
- [5] Lange, O., "The Determinateness of the Utility Function," *The Review of Economic Studies*, I, 218-225 (1933-34).
- [6] Neumann, J. von and O. Morgenstern, Theory of Games and Economic Behavior, Science Edition (John Wiley and Sons, Inc., New York, N. Y., 1964).
- [7] Pareto, V., Manuel d'économie politique (Paris, 1927).
- [8] Ramsey, F. P., The Foundations of Mathematics (The Humanities Press, New York, 1950).

\* \* \*





# AN EVALUATION OF INCENTIVE CONTRACTING EXPERIENCE

I. N. Fisher

*The RAND Corporation*

## ABSTRACT

Incentive contracts are intended to motivate defense contractors to perform more efficiently and control costs more closely. By increasing the total profit as actual costs are reduced below a predetermined cost target, they encourage contractors to achieve cost underruns. Consequently, the principal advantage claimed for these contracts is that they make the financial incentives to reduce costs more effective.

This study examines the effectiveness of incentive contracts as a means for controlling defense procurement costs. The study considers the various effects that incentive contracts may have on both contractors' performance and contract costs, and presents empirical evidence suggesting that these contracts may not have accomplished their intended goal of increased efficiency and lower procurement costs.

## I. INTRODUCTION

Present defense procurement policy relies heavily on incentive pricing arrangements to motivate contractors to control costs and perform more efficiently — increased profits being the incentive.

There seems to be general agreement among Defense Department officials that incentive contracts have achieved their goal. Nonetheless, there are reasons for questioning the cost-saving effects claimed for these contracts. This study examines the effects of incentive contracts on contract costs and considers some prospects for improving their effectiveness.

Two basic types of pricing arrangements are used in defense contracting: fixed-price and cost-reimbursable contracts. Under cost-reimbursable contracts, the Government pays all legally allowable costs the contractor incurs during the life of the contract. Under fixed-price agreements, the contractor and the Government agree on a target cost; any discrepancy between the target cost and the actual cost may be borne entirely by the contractor or shared in some predetermined proportion with the Government.

The following are the major varieties of pricing arrangements:

### Fixed-price contracts

- Firm-fixed-price (FFP)
- Fixed-price-incentive (FPI)
- Fixed-price-redeterminable (FPR)

### Cost-reimbursable contracts

- Cost-plus-fixed-fee (CPFF)
- Cost-plus-incentive-fee (CPIF)

Most fixed-price contracts are presently either fixed-price-incentive (FPI) or firm-fixed-price (FFP).<sup>\*</sup> FFP contracts presumably provide the maximum incentive for contractors to control cost, since the price remains fixed once the target has been established. FPI contracts, in contrast, contain a profit-sharing arrangement whereby the Government and contractor share any difference that occurs between the actual and target cost.

Only two varieties of cost-reimbursable contracts are used extensively. The cost-plus-fixed-fee (CPFF) contract reimburses the contractor for all allowable costs incurred in completing the contract. The contractor also receives a fixed fee that does not depend on his cost performance. Obviously, this form of pricing arrangement provides little, if any, incentive for contractors to control costs.

The cost-plus-incentive-fee contract (CPIF) is also a cost-reimbursable contract, since the Government bears all the allowable costs of contract performance; it is also an incentive contract, since it establishes both a target cost and a profit-sharing arrangement.

In 1962, the Defense Department revised the Armed Services Procurement Regulations (ASPR) to encourage the use of incentive contracts. These revisions reflected a consensus within the Defense Department that the CPFF contracts then commonly used to purchase major weapon systems did not provide adequate incentive for contractors to control costs. These revisions established CPIF contracts as preferable for research and development effort and recommended the use of FFP or FPI contracts for production. Use of CPFF contracts is limited to situations involving considerable uncertainty, in which incentive-type contracts would be impractical.

These changes have had a tremendous impact on the defense industry and have resulted in a substantial increase in the use of FFP and other types of incentive contracts for defense procurement. As Table 1 indicates, the shift away from CPFF contracts has been striking. CPFF contracts accounted for more than one-third of total defense expenditures in 1960, but less than 10 percent in 1966. In the same period, CPIF contracts more than doubled, and FFP contracts nearly doubled.

Defense Department representatives believe incentive contracts are effective in controlling procurement costs. The principal advantage claimed for these contracts is that they make the financial incentives to reduce costs more effective. By increasing the total profit as actual costs are reduced below the target, they encourage contractors to achieve cost under-runs. They also place greater financial risk on the contractor since the Government no longer stands ready to absorb cost overruns completely.

Cost overruns have been far less frequent and less substantial under incentive contracts than under CPFF contracts. Defense Department officials have interpreted this outcome as evidence that a contractor's performance under incentive contracts is more efficient than under CPFF contracts. In fact, in evaluating the impact of incentive contracts on procurement costs, former Secretary McNamara stated that costs under incentive contracts are ten percent lower than they would have been under CPFF pricing arrangements.<sup>†</sup> Nonetheless, there are some valid reasons for questioning the extent of the cost savings claimed for these

<sup>\*</sup>Fixed-price-redeterminable contracts (FPR) are no longer extensively used; moreover, it is erroneous to classify them as fixed-price contracts, since they provide for periodic price renegotiation during the life of the contract.

<sup>†</sup>See Statement of Secretary of Defense Robert S. McNamara Before the House Armed Services Committee on the Fiscal Year 1966-1970 Defense Program and 1966 Defense Budget, February 18, 1965, Senate Subcommittee on DOD Appropriations, p. 187.

TABLE 1<sup>a</sup>  
Percentages of Total Defense Expenditures by  
Type of Pricing Arrangement

Contract Type	Fiscal Year						
	1960	1961	1962	1963	1964	1965	1966
Fixed Price							
FFP	31.4	31.5	38.0	41.5	46.3	52.8	57.5
FPI	13.6	11.2	12.0	15.8	18.5	16.6	15.9
Other <sup>b</sup>	12.4	15.2	10.8	7.6	6.4	7.1	5.8
Cost-Reimbursable							
CPFF	36.8	36.6	32.5	20.7	12.0	9.4	9.9
CPIF	3.2	3.2	4.1	11.7	14.1	11.2	8.3
Other <sup>c</sup>	2.6	2.3	2.6	2.7	2.7	2.9	2.6

<sup>a</sup>Source: Directorate for Statistical Services, OSD, Military Prime Contract Awards.

<sup>b</sup>Includes FPR contracts.

<sup>c</sup>Includes cost and cost-sharing contracts.

contracts. The most important reason is that cost underruns often may be achieved without any real cost savings to the Government.

## II. STRUCTURE OF INCENTIVE PRICING ARRANGEMENTS

Incentive contracts are supposed to motivate defense contractors to perform more efficiently and control costs more closely. This is accomplished through the incentive sharing provision, which allows contractors to retain part of any resulting cost underrun as increased profits. So long as these underruns represent realized cost reductions, incentive contracts accomplish their intended goal.

To understand how cost underruns may occur without benefit of real cost savings to the Government, consider the factors that determine the contractor's profit under an incentive contract. Total profit received by the contractor consists of two components,

$$\Pi_T = \Pi_t + \alpha(C_t - C_f) ,$$

where

$\Pi_T$  = total fee to contractor,

$\Pi_t$  = profit on initial target amount,

$C_t$  = target cost,

$C_f$  = actual cost, and

$\alpha$  = incentive sharing rate.



The first component is the profit amount based on the target cost. The second component is the profit sharing arrangement by which contractors retain part of any cost underrun that may result, but must bear a portion of any cost overrun. The term inside the parentheses is an overrun when the actual cost exceeds the target, and an underrun when actual cost is less than the target.

The incentive feature operates through this profit-sharing arrangement. To obtain increased profits, the contractor must achieve a cost underrun. For each dollar increase in underrun, the contractor retains  $\alpha$  percent as increased profit, providing motivation to achieve as large an underrun as possible.

It is of course in the contractor's interest to increase the underrun and thereby increase his profit. One way he does this is to perform more efficiently and hold actual costs below the target value — the effect desired by Defense Department officials. Because overruns and underruns depend on both the actual cost and the target cost, another method for avoiding overruns and increasing underruns is to secure as high a target cost as possible. The success of this strategy of course depends on the circumstances under which the target is determined. So long as targets are determined competitively, procurement officials need have little concern over their precise values. The market forces operating in a competitive environment tend to nullify the possibility of obtaining targets that are in some sense excessive.\*

The difficulty is in determining appropriate target values for contracts negotiated in a noncompetitive environment. This problem is significant, because most weapon system production and support contracts are presently negotiated in the absence of any price competition. Moreover, many development contracts that are awarded competitively are awarded on the basis of technical or nonprice rivalry. Because target costs are commonly negotiated in these situations, contractors have much greater opportunity to increase them. If they succeed, the resultant targets may fail to provide any real incentives for cost reduction and efficiency.

Provided the Government has adequate information upon which to predict cost as well as the technical expertise required to make an independent cost estimate, a realistic target can be negotiated. Otherwise, an inflated target and a consequent underrun are the likely results. Such an underrun is unrelated to any real cost savings, merely reflecting the larger target cost.†

\*Although competition may eliminate excessive target costs, it may also result in the selection of a less efficient contractor. See J. J. McCall, An Analysis of Military Procurement Policies, The RAND Corporation, RM-4062-PR, November 1964.

†The supplemental changes and modifications that occur after the target has been established also provide an opportunity for the contractor to increase the target cost above the expected value. More precisely, the profit formula should be written

$$\Pi_T = \Pi_t + \Pi_S + \alpha(C_a - C_f),$$

where

$\Pi_S$  = additional fee allowed on supplemental changes and modifications;

$C_a$  = adjusted target cost, including the negotiated costs of supplemental changes and modifications.

It is apparent that incentive pricing arrangements may also encourage contractors to propose frequent changes and modifications to the initial contract because these changes may result in additional profits,  $\Pi_S$ . Moreover, since the costs of changes and modifications must be negotiated, it also provides an opportunity for the contractor to increase the target cost, thereby improving the likelihood of an underrun. Through the remainder of this section, the term "target cost" will include the effect of supplemental changes and modifications, i.e., the adjusted target cost.



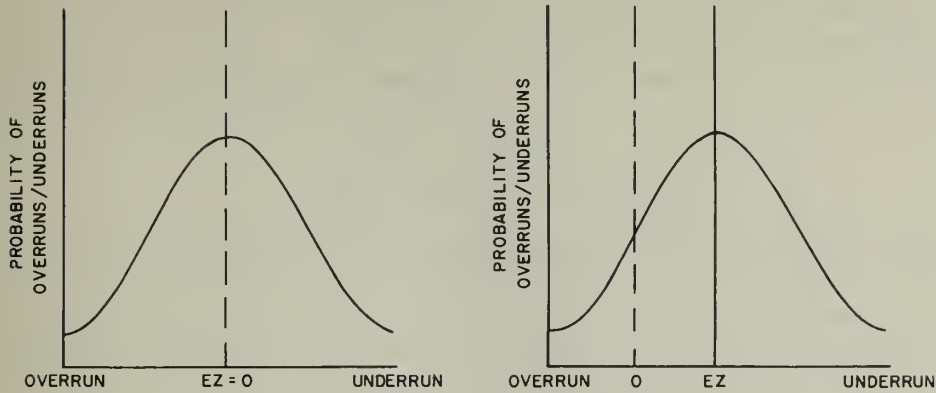


Figure 1. Distribution of cost overruns/underruns

This point is illustrated in Figure 1. Assume that actual cost,  $C_f$ , is a random variable with probability distribution  $\phi(EC_f, \sigma^2)$ , where  $EC_f$  is the expected actual cost and  $\sigma^2$  is the variance. An overrun/underrun is defined as  $Z = C_t - C_f$ , where  $C_t$  is the target cost, which is a constant. Now  $Z$  is a random variable also with probability distribution given by  $\phi(EZ, \sigma^2)$ , where  $EZ$  is the expected overrun/underrun and is equal to  $(C_t - EC_f)$ . In (a), the negotiated target cost is equal to the expected actual cost so that the expected overrun/underrun,  $EZ$ , is zero. That is, overruns and underruns are likely to occur with equal probability (at least for probability distributions symmetrical about  $EZ$ ). In (b), however, the target cost is larger than the expected final cost so that  $(C_t - EC_f) > 0$ . In this case, underruns are more likely than overruns (again, for symmetrical distributions).

There is a second reason why target costs may be larger with incentive pricing arrangements. Incentive contracts increase the risk of financial loss to the contractor by requiring him to bear part of any cost overrun that may result. Assuming that contractors are generally averse to risk, profits on incentive contracts must be sufficient to offset the increased risk.

This is illustrated in Figure 2. Curves  $U_1$  and  $U_2$  are typical indifference curves for a risk-averse contractor. These curves are upward-sloping, indicating that larger profits are required to compensate for increased financial risk, measured in terms of the sharing rate value.  $U_2$  represents greater utility or satisfaction than  $U_1$ ; given a contract with sharing rate  $\alpha_1$ , for example, an increase in expected profit from  $\Pi_1$  to  $\Pi_2$  increases the contractor's utility from  $U_1$  to  $U_2$ .

Consider a minimum-risk CPFF contract with expected profit equal to  $\Pi_1$ . Since the sharing rate value for CPFF contracts is zero, the level of utility corresponding to this contract is  $U_2$ . Now suppose the Government replaces it with an incentive contract having a sharing rate value equal to  $\alpha_1$ . If the expected profit is unchanged, the contractor's utility decreases from  $U_2$  to  $U_1$ . In order to make the incentive contract equally attractive to the contractor, expected profit must be increased from  $\Pi_1$  to  $\Pi_2$ . At this profit level the contractor will be indifferent between the two contracts. The profit differential,  $\Pi_2 - \Pi_1$ , is the

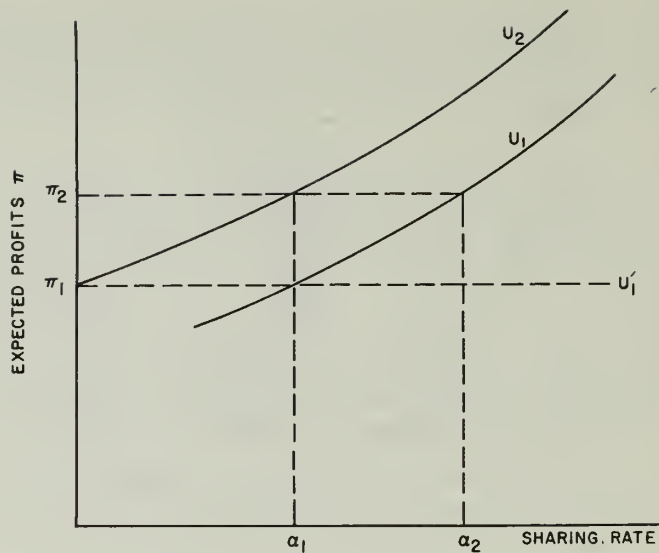


Figure 2. Contractor's indifference curves

risk premium necessary to compensate the contractor for the increased financial risk attached to the incentive contract.\*

Compensation for the increased risk attached to incentive contracts can be provided in several ways. An obvious method would be for the Government to increase target profits by the amount of the required risk differential.† In practice, however, it may not be possible to increase profits sufficiently to offset the increased risk completely. If risk is large and the contractor extremely risk-averse, the required risk premium may be so large as to result in a rate of profit that is politically prohibitive. For example, profit rates of 40 percent or more might be required on some contracts; such rates would arouse Congressional interest and be difficult to explain.

Since it may not be possible to increase profits sufficiently to offset the increased financial risk inherent in incentive contracts, contractors may reduce the risk level by negotiating target costs high enough to provide a margin of safety against large overruns. This strategy is justified whenever profits are not sufficient to offset the greater risk completely.‡ In short, both larger profits and larger target costs may be required to compensate for the greater risk attached to incentive contracts. Unfortunately, the underruns that accompany these larger targets may be erroneously attributed to reduced costs and increased efficiency.

\*If the contractors were neutral toward increased risk, the utility function would appear as curve  $U_1'$ . In this situation contractors would be indifferent between minimum-risk CPMF contracts and maximum-risk FFP contracts, and no risk premium would be required.

†The DOD has recognized the need for larger profits on riskier incentive contracts with larger sharing rate values; the ASPR specifically provides for larger negotiated profit rates for these contracts. See ASPR 3-808.1 (6).

‡Evidence indicates that larger target costs are negotiated as the sharing rate becomes larger. See John Cross, "A Reappraisal of Cost Incentives in Defense Contracting," P-282, Institute for Defense Analyses, 1966, and F. M. Scherer, *The Weapons Acquisition Process: Economic Incentives* (Harvard University Press, Boston, 1963). This has been explained as the compensation required to induce contractors to bear greater risk. Nonetheless, larger targets reduce the probability of overruns and increase the likelihood of increased profits.

This discussion brings out the important point that incentive contracts really provide two different incentives; not only do they motivate contractors to reduce actual costs, but they also encourage them to overstate target cost estimates.\* Thus, it may be misleading to attribute the underruns observed with these contracts to reduced costs and improved performance without more detailed analysis of the available evidence.

### III. SOME EMPIRICAL EVIDENCE

One limitation common to all empirical analyses of incentive contracts is that it is never clear whether the underruns observed with incentive contracts result from increased efficiency and better cost control or from larger targets secured by contractors to compensate for increased risk. To assess the true impact of incentive pricing arrangements on the cost of military procurement, it would be necessary to determine how incentive pricing provisions affect target costs. The data required for the analysis, however, are not available. Nonetheless, although it is impossible to separate directly the effects that incentive contracts have on target costs from their effects on actual costs, it is possible to draw some inferences about how contractors respond to these contracts by examining several other measures of cost outcome for which data are available.†

Three questions are examined in this section. These are: (1) whether the underruns observed with incentive contracts are related to the incentive features of these contracts; (2) whether these underruns result from supplemental changes and modifications that occur during the life of the contract; and (3) the extent to which observed underruns differ among major defense contractors. The statistical analyses presented here are somewhat technical; however, a summary and discussion of the major conclusions and implications are presented at the end of the section.

#### COST OVERRUNS/UNDERRUNS

Since incentive pricing arrangements sharpen the incentive for contractors to seek cost underruns, one would naturally expect to find that underruns are more common with incentive contracts than with cost-reimbursable contracts. Table 2 compares the average cost overrun/underrun for several types of contracts included in the sample. An average overrun is observed

\*The relative importance of these two incentive effects depends on the values of the incentive sharing rate and the rate of profit allowed on the contract. For example, differentiating the profit function with respect to both target cost and actual cost yields:

$$dP/dC_t = (p_t + \alpha)$$

and

$$dP/dC_a = -\alpha.$$

The first term is the marginal effect on profits from a change in the target cost; the second is the marginal effect of a change in the actual cost. Since  $dP/dC_t > 0$ , an increase in the target cost results in an increase in the total profit. On the other hand, since  $dP/dC_a < 0$ , an increase in actual cost reduces the total profit. Since  $(p_t + \alpha) \geq \alpha$ , the effect of increasing the target cost by one dollar outweighs the effect of reducing actual costs by the same amount, and as long as  $p_t > 0$ , the incentive to overstate target costs will be more tempting than will be the incentive to reduce actual costs.

†The sample used in this section contains 1007 Air Force contracts completed during fiscal years 1959 through 1966. The data consist solely of contracts for major weapon systems and related equipment and total nearly \$15.7 billion. For a more detailed description of the sample characteristics, see I. N. Fisher, A Reappraisal of Incentive Contracting Experience, The RAND Corporation, RM-5700-PR, July 1968.

TABLE 2  
Average Overrun/Underrun by Type of Contract<sup>a</sup>  
Percentage of Final Cost

Type of Contract	Mean Overrun/Underrun
FPI	-3.18
FPR	1.74
CPIF	1.29
CPFF	1.90

<sup>a</sup>Unweighted averages of observed overruns/under-runs for each type of pricing arrangement.

for each group except the fixed-price-incentive contracts (FPI), illustrating the trend that is often interpreted by Defense Department officials as an indication of greater efficiency and cost reduction.\*

#### Pricing Arrangement

Table 2 indicates that average overruns/underruns are different for the FPI as opposed to the other three groups. The significance of these differences can be tested statistically using analysis of variance to determine whether the observed overruns/underruns differ significantly among the two groups of contracts.<sup>†</sup> Table 3 presents the results. Note that the mean-square deviation of overruns/underruns is large between groups and small within groups. This indicates that there are significant differences in overrun/underruns between the two groups, but little variation within each group. An F-ratio value greater than 3.78 (at the 0.01 level of probability) is required in order to be confident that the observed differences among the two groups are anything but spurious. Since the computed value of the F-ratio is 18.5, the analysis indicates that these observed differences are statistically significant and are unlikely to have occurred by chance.

#### Incentive Sharing Rate

Table 4 shows the average overrun/underrun and its standard deviation for incentive contracts classified according to sharing rate value. Note that an average overrun occurs for contracts in the first group-- those with sharing rate values less than 10 percent -- while the remaining three groups have average underruns. The reason is that most of the contracts in the first group are CPIF, while those in the remaining groups are FPI and FPR and, on the basis of the preceding results, overruns would be expected on average for the CPIF group. Although average underruns are indicated for the three groups with sharing rate values greater

\*Overrun/underrun is computed from  $(C_f - C_a)/C_f$ , so that underruns are negative while overruns are positive.

<sup>†</sup>FFP contracts were excluded since no measure of overrun/underrun is available for these contracts. For FPR, the price is periodically renegotiated during the life of the contract. As a result, these arrangements closely resemble cost-reimbursable contracts. This is also indicated to some extent in Table 2, where the average overrun observed for FPR contracts is nearly as large as that for CPFF.



TABLE 3  
Analysis of Variance for Two Contract Classifications:  
FPI and All Others<sup>a</sup>

Variance	Sum of Squares	D. F.	Mean Squares	F Ratio
Between group	0.3826	1	0.3826	18.573
Within group	19.4783	946	0.0206	
Total	19.8609	947		$F_{0.01} = 6.64$

<sup>a</sup>In order to determine whether the observed overrun/underruns are statistically different between groups, the within-group variation (i.e., variation of overruns/underruns in each group about the group mean) is compared with the between-group variation. If the variation within the groups is large while that between groups is small, differences between the groups may be insignificant. On the other hand, small within-group variation but large between-group variation suggests that the observed differences between groups may be significant. Analysis of variance computes the ratio of these variations (adjusted for degrees of freedom) and provides a formal method for testing the significance of the ratio.

TABLE 4  
Mean Overrun/Underrun and Standard Deviation  
CPFF and FFP Contracts Excluded

Item	Sharing Rate Value			
	0.01-0.09	0.10-0.19	0.20-0.29	0.30-0.99
Mean <sup>a</sup>	1.45	-3.50	-2.32	-0.39
Standard deviation <sup>a</sup>	12.95	13.86	8.45	8.81
Number	43	144	156	87

<sup>a</sup>Measured as a percentage of final cost.

than 0.10, they appear to become progressively smaller as the sharing rate becomes larger. This is a curious result, since the opposite trend would be expected; i.e., the larger sharing rates presumably subject the contractor to greater risk of financial loss and, consequently, provide stronger motivation to avoid overruns. The large standard deviations for each of these groups, however, indicate that there is considerable variation about the mean values.

Figure 3 illustrates one possible way of describing the predicted relationship between the sharing rate and cost overruns/underruns. If conventional beliefs about incentive fees are correct, low sharing rate values should be associated with cost overruns, while larger sharing rate values should be associated with cost underruns. This type of linear relationship can be described by an equation of the form:



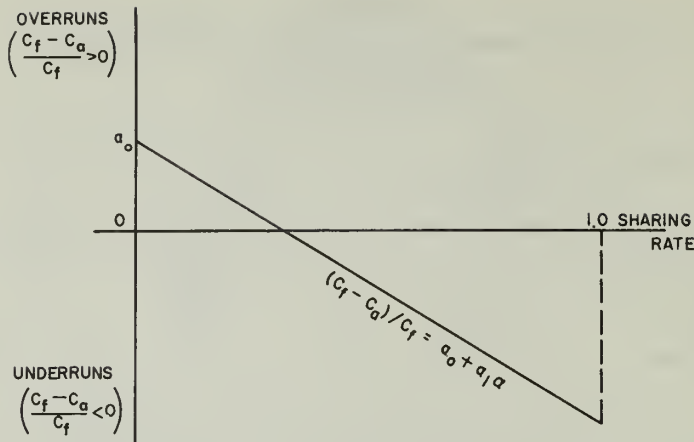


Figure 3. Relationship between incentive sharing rate and overruns/underruns

$$(1) \quad (C_f - C_a)/C_f = a_0 + a_1 \alpha,$$

where

$\alpha$  = incentive sharing rate,

$C_f$  = final cost,

$C_a$  = adjusted target cost (initial target plus changes and modifications), and

$a_0, a_1$  are undetermined coefficients.

One procedure for testing whether underruns are, in fact, larger for contracts with larger sharing rate values would be to compute the correlation between these two variables. A more interesting procedure is to estimate the values of the coefficients,  $a_0$  and  $a_1$ , in Eq. (1) by means of simple regression analysis. This provides a measure of correlation along with several other useful statistics.

This has been done for the FPI, FPR, CPIF contracts; the results appear in Table 5. The value of the constant term,  $a_0$  (shown in column 1), is the intercept illustrated in Fig. 3. These values provide an estimate of the average overrun/underrun that would result for the three types of contracts if the sharing rate were to approach zero. The estimated value for  $a_1$ , shown in column 2, is the slope of the curve illustrated in Fig. 3, and provides a measure of the effect of the sharing rate value on cost overruns/underruns. If underruns are larger for larger sharing rate values, the coefficient  $a_1$  should be negative. A positive sign for this coefficient would imply smaller underruns for contracts with larger sharing rate values. The numbers appearing in the third column are the standard errors of the slope estimates and provide a measure of their reliability. The last column contains a measure of the correlation between observed overruns/underruns and the sharing rate value.

Although two of the three estimated values for  $a_1$  shown in Table 5 are negative, the coefficient for the FPI contracts is positive. Moreover, when these three estimates are compared with their standard errors, shown in column 3, the FPI coefficient is the only one that is

TABLE 5  
Estimated Regression Coefficients  
Cost Underrun/Overrun and Sharing Rate

Type of Contract	$a_0$	$a_1$	Standard Error of $a_1$	$R^2$
FPI	-0.100	0.3167	0.1066	0.0270
FPR	0.084	-0.7501	0.4441	0.0754
CPIF	0.050	-0.1903	0.1313	0.0283

statistically significant (at the 0.01 level of probability). This surprising result suggests that underruns are smaller, not larger, for FPI contracts with larger sharing rate values.\*

#### Contract Size

Another important factor that may affect this relationship is contract size. It seems reasonable to expect contractors to be more concerned with multimillion-dollar contracts than with million-dollar contracts because the financial consequences of cost overruns are much more serious. This possibility can be investigated by including a measure of contract size in Eq. (1). One obvious measure that could be used is contract cost. Including this in the relationship results in

$$(2) \quad (C_f - C_a)/C_f = a_0 + a_1\alpha + a_2C_f',$$

where  $C_f' = \log$  of final cost,  $C_f$ .<sup>†</sup>

Estimates of these coefficients for each type of contract appear in Table 6. None of the coefficients of size,  $a_2$ , are statistically significant at any reasonable level of confidence. Moreover, including size in the relationship has had little effect on the value of  $R^2$  (compare Table 5). Consequently, size appears to have little effect on observed overruns/underruns for incentive contracts.

TABLE 6  
Estimated Regression Coefficients  
Cost Overrun/Underrun, Sharing Rate, and Contract Size

Type of Contract	$a_0$	$a_1$	Std. Error of $a_1$	$a_2$	Std. Error of $a_2$	$R^2$
FPI	-0.174	0.311	0.107	0.019	0.016	0.032
FPR	0.128	-0.166	0.133	-0.021	0.032	0.093
CPIF	-0.019	-0.713	0.449	0.026	0.031	0.045

\*This result holds only for the sample of Air Force contracts examined here. Nonetheless, it is based on a fairly large number of contracts and is difficult to explain. This question will be examined in more detail, however, in a later section discussing the observed pattern of supplemental changes.

<sup>†</sup>Since costs vary widely (between one million and several hundred million), the logarithm of final cost has been used in place of the final cost. This acts as a scale factor to reduce large absolute differences in dollar amount while preserving relative differences.

### Summary

These results indicate that although underruns are more common for FPI contracts than for other types, the underruns do not seem to be related to either the value of the sharing rate or to the size of the contract. Cost overruns/underruns appear to be no different for contracts with small sharing rate values than for those with large sharing rates, or for contracts differing substantially in total dollar amount. This suggests that these contracts have not had an important effect on contract cost outcomes or contractor performance.

Since the magnitude of the cost overruns/underruns observed with FPI contracts seems to be unrelated to either the value of the sharing rate or to contract size, it is difficult to attribute these underruns to increased efficiency and reduced costs. It is unlikely that contractors are equally efficient and cost-conscious for all FPI contracts regardless of differences in financial risk associated with the sharing rate and size of contract. It is more likely that these observed underruns result from larger target costs — targets that exceed anticipated actual costs.

There are essentially two ways in which contractors could insure that the adjusted target cost exceeds the expected cost. One would be to negotiate larger target costs to the extent possible during contract negotiation. This, of course, would depend on the circumstances under which the contract was awarded — that is, the degree of monopoly power enjoyed by the contractor. Another possibility, once the initial target had been negotiated, would be to introduce numerous and costly changes and modifications in the original specifications. This strategy would also improve the likelihood of achieving cost underruns. Although there is presently no way of determining how inflated initial target costs may be, the costs of supplemental changes and modifications are known. Consequently, the effects of different contract characteristics on the magnitude of these costs can be explored in some detail.

### SUPPLEMENTAL CHANGES AND COST UNDERRUNS

Table 7 summarizes average costs of supplemental changes, measured as a percentage of final cost, for four major types of contracts. Supplemental changes appear to be considerably larger for the cost-reimbursable contracts than for the fixed price contracts. This may reflect the greater technical uncertainty inherent in those projects typically included under CPFF coverage.

TABLE 7  
Summary Statistics: Supplemental  
Changes and Modifications  
Percentage of Final Cost

Type of Contract	Mean
FPI	4.17
FPR	7.97
CPIF	77.15
CPFF	60.08

Contractors may be able to increase the likelihood of achieving a cost underrun by introducing frequent supplemental changes and modifications. This strategy provides the opportunity to adjust the target cost upward and would appear particularly attractive whenever the target cost is tight; i.e., close to the contractor's anticipated actual cost. One way of investigating this possibility is to estimate the relationship given in Eq. (3):

$$(3) \quad (C_f - C_a)/C_f = a_0 + a_1(C_a - C_i)/C_f,$$

where

$C_f$  = final contract amount,

$C_a$  = adjusted target cost (initial target plus supplemental changes),

$C_i$  = initial negotiated target cost,

and  $a_0$ ,  $a_1$  are unknown coefficients to be estimated. If underruns are greater for contracts with larger supplemental changes,  $a_1$  should be negative and statistically significant.

Table 8 presents the estimated values for the coefficients along with their standard errors and coefficients of determination for three types of incentive contracts. In all cases the values of  $R^2$  are extremely low, indicating that underruns and supplemental changes are not closely related. Thus it appears that contractors do not utilize supplemental changes to inflate target costs and increase the magnitude of cost underruns.

TABLE 8  
Estimated Regression Coefficients: Underruns and Supplemental Changes

Type of Contract	$a_0$	$a_1$	Standard Error of $a_1$	$R^2$
FPI	-0.0274	0.01002	0.03093	0.0003
FPR	0.0085	-0.17541	0.08393	0.0571
CPIF	-0.0164	-0.00862	0.04072	0.0012

#### COMPARISON OF OVERRUNS/UNDERRUNS AMONG CONTRACTORS

The results obtained above indicate that both cost overruns/underruns and supplemental changes differ markedly between cost-reimbursable and fixed-price contracts. Thus, it appears that contractors react differently to these two types of pricing arrangements. There may also be significant differences in cost performance among individual contractors, however. For example, some contractors may achieve cost underruns consistently, while overruns may be typical for others.

Table 9 compares the average overrun/underrun for those contracts exceeding \$1 million held by 15 large Air Force contractors. As before, underruns are more common for the fixed-price contracts (FPI and FPR) than for the cost-reimbursable contracts (CPIF and CPFF). Nonetheless, several contractors have average underruns for both cost-reimbursable and fixed-price contracts, while others have average overruns for both types. This suggests that there may be some important differences among individual contractors' responses to these contracts.



TABLE 9  
Comparison of Average Overruns/Underruns: Fifteen Large Air Force Contractors

Contractor	Cost-Reimbursable Contracts		Fixed-Price Contracts	
	Average Overrun	No.	Average Overrun	No.
1	0.0463	10	0.0323	13
2	-0.0436	20	-0.0746	21
3	0.0671	7	0.0010	10
4	-0.0044	21	-0.0106	13
5	0.0142	29	-0.0327	18
6	0.0012	51	-0.0507	4
7	-0.0412	11	0.0107	72
8	-0.0080	50	0.0598	6
9	0.0050	13	-0.0097	15
10	-0.0889	28	-0.0219	29
11	-0.0958	21	-0.0726	6
12	0.1267	4	-0.0220	31
13	0.0325	40	-0.0950	11
14	0.0631	10	-0.1595	6
15	0.0086	10	-0.0421	32

Analysis of variance can be used to determine whether these average overruns/underruns differ significantly between contractors. Tables 10 and 11 present the results for cost-reimbursable and fixed-price contracts, respectively. For the cost-reimbursable contracts, the differences in average overruns/underruns between contractors are not statistically significant. That is, overruns/underruns are not noticeably different among cost-reimbursable contracts for any contractors included in the sample. For the fixed-price incentive contracts, however, the value of the F ratio is statistically significant (at the 0.01 level of probability), indicating that there are important differences in observed overruns/underruns among contractors. Thus, for fixed-price contracts, some contractors apparently experience larger cost underruns, on an average, than do others.

There are two possible explanations for these differences. It may be that some contractors are more responsive to contract profit incentives than others; i.e., some contractors may perform more efficiently or apply greater pressure for larger target costs on contracts with larger sharing rates. On the other hand, some contractors may be generally more efficient than others regardless of the pricing arrangement, may be more aggressive in negotiating larger costs, or may enjoy certain competitive advantages that increase the likelihood of achieving cost underruns. For example, both an absence of market price information and lack



TABLE 10  
Analysis of Variance: Cost-Reimbursable Contracts,  
Fifteen Large Air Force Contractors

Variance	Sum of Squares	D. F.	Mean Squares	F
Between group	0.5597	14	0.0400	1.596
Within group	7.7668	310	0.0251	
Total	8.3265	324		$F_{0.05} = 1.72$

TABLE 11  
Analysis of Variance: Fixed-Price Contracts,  
Fifteen Large Air Force Contractors

Variance	Sum of Squares	D. F.	Mean Squares	F
Between group	0.4915	14	0.0351	4.228
Within group	2.2596	272	0.0083	
Total	2.7511	286		$F_{0.01} = 1.73$

of meaningful competition improve the contractors' ability to obtain larger target costs and larger cost underruns.

If the first explanation -- differences in contractors' responses to incentive pricing arrangements -- is correct, observed underruns for each contractor should be larger for contracts with stronger profit incentives. Alternatively, if the observed differences in cost underruns result principally from negotiation strategy, market position, and overall efficiency, the magnitude of the underruns should be relatively constant for a given contractor.

One way to determine which alternative better explains the observed differences in underruns among contractors is to estimate the relationship between the pricing arrangement and cost overrun/underrun for each contractor. This relationship is described in Eq. (4).

$$(4) \quad (C_f - C_a)/C_f = \gamma_j + b_j \alpha \quad j = 1, \dots, 15,$$

where

$\gamma_j$  = average overrun/underrun for the  $j$ th contractor, and

$b_j$  = effect of the incentive sharing rate on the  $j$ th contractor.

In this formulation the overrun/underrun for each contract, measured as a percentage of final cost, is expressed as the sum of two components; the first is the average overrun/underrun for the individual contractor,  $\gamma_j$ , while the second,  $b_j$ , reflects the effect

of the pricing arrangement on the contractor. If the estimated values for the  $b_j$ 's differ significantly among contractors while the values for the  $\gamma_j$ 's remain relatively constant, then variations in observed underruns should be attributed to differences in individual contractors' responses to the profit incentives. If, on the other hand, the  $b_j$ 's are relatively constant but the  $\gamma_j$ 's vary noticeably among contractors, then variations in observed overruns/underruns should be attributed to individual contractor-specific characteristics such as the contractor's ability to estimate target costs, his competitive advantage, negotiation strategy, and overall differences in efficiency.

Estimates of these coefficients appear in Table 12. Note the differences between the average overruns/underruns shown in Table 9, and the estimated values shown here. These differences occur because the sharing rate accounts for a portion of the average overruns/underruns shown in Table 9. The estimated values for both the  $\gamma_j$  and  $b_j$  coefficients differ substantially among contractors. The significance of these variations can be determined by testing the following hypotheses:

$$H_1: \gamma_i - \gamma_j = 0 \quad \text{for all } i, j;$$

$$H_2: b_i - b_j = 0 \quad \text{for all } i, j.$$

Hypothesis 1 asserts that there are no significant differences in average overruns/underruns among contractors, while Hypothesis 2 asserts that the effect of the pricing arrangement on overruns/underruns is negligible for each contractor.

TABLE 12  
Estimated Coefficients: Fixed-Price Contracts

Contractor	$\gamma_j$	$b_j$	Contractor	$\gamma_j$	$b_j$
1	0.0836	0.0081	9	-0.0179	0.0046
2	-0.1541	-0.0159	10	-0.0175	0.0059
3	0.0313	0.0055	11	-0.0279	-0.0099
4	0.0234	-0.0029	12	-0.0712	0.0033
5	-0.0117	-0.0056	13	-0.0945	0.0153
6	-0.0285	0.0041	14	-0.1242	0.0080
7	0.0294	0.0011	15	-0.0518	-0.0158
8	-0.0933	-0.0122			

These hypotheses are also tested using analysis of variance; the ratios of the explained mean square deviation to the unexplained mean square deviation for each set of variables are computed in Table 13. Since the critical value of the F ratio (at the 0.01 level) is 1.79,  $H_1$  can be rejected while  $H_2$  cannot. This means that the  $b_j$ 's are not statistically different from zero or, in other words, that the incentive pricing arrangement has had little effect on the cost performance of the contractors. Thus the overruns/underruns observed for these contractors must be explained by other factors.

TABLE 13  
Analysis of Variance: Fixed-Price Contracts

Item	Sum of Squares	D. F.	Mean Squares	F
Pricing arrangement, $b_j$	0.0754	15	0.0050	0.549
Contractor effect, $\gamma_j$	0.3347	15	0.0223	2.451
Unexplained residual	2.3411	256	0.0091	
Total variation	2.7512	286		

The  $\gamma_j$ 's, on the other hand, are statistically significant, indicating that there are important differences in average overruns/underruns among the contractors. Thus, some contractors consistently achieve larger underruns than others and, since these underruns cannot be explained by differences in incentive pricing arrangements, they must be related to other characteristics peculiar to each contractor.

These results indicate that although average overruns/underruns for fixed-price contracts differ significantly among contractors, these differences cannot be explained by variations in contract pricing arrangements. There appears to be no relationship between the incentive arrangement and the observed cost overruns/underruns for any of the individual contractors examined. Consequently, it seems improbable that the larger underruns achieved by some contractors result from increased efficiency or reduced costs.

It may be that contractors with the largest underruns produce different types of products involving less uncertainty than do those experiencing smaller underruns (or larger overruns). This explanation seems unlikely, however, since all 15 contractors examined were large, well diversified, and produced similar products. It seems more likely that these observed underruns result from differences in other contractor-related factors — factors that include competitive advantage, cost-estimating ability, negotiation skill, and general managerial capability. Some contractors may consistently be able to obtain larger target costs than others, for example, thereby increasing the likelihood of obtaining underruns.

## CONCLUSIONS

In sum, these results indicate that although underruns are more common with fixed-price-incentive contracts, they are not related to the pricing provisions of the contract. Consequently, these underruns should not be attributed to increased efficiency or reduced costs. It is difficult to believe that contractors are generally more efficient and cost-conscious under FPI contracts regardless of differences in financial risk associated with the incentive sharing rate. It seems more likely that these observed underruns result primarily from target costs that exceed anticipated actual costs.

Contractors could increase the adjusted target cost above the anticipated actual cost by either of two possible strategies. One would be to negotiate larger initial target costs — targets that are sufficiently greater than expected actual costs to provide a margin of safety against possible cost overruns. The extent to which this may be possible, of course, depends on the degree of price rivalry as well as on the Government's ability to predict actual costs accurately.

The other alternative would be to introduce numerous supplemental changes in order to provide a basis for negotiating a larger adjusted target cost, thereby increasing the likelihood of an underrun. The results obtained here, however, indicate that supplemental changes do not explain the underruns observed with incentive contracts; these underruns seem to be generally unrelated to the magnitude of the supplemental changes. The evidence indicates that these observed underruns originate, instead, from target costs that exceed the contractor's anticipated actual cost. Given present weapon system procurement practices, it is easy to see how this may occur. So long as subsequent production and follow-on contracts are awarded to the initial development contractor without effective price rivalry, there can be no guarantee that the negotiated target cost is sufficiently close to the contractor's anticipated actual cost to provide a meaningful incentive for greater efficiency and reduced costs.

In short, incentive contracts cannot be expected to provide the motivation for which they were intended without some means for establishing realistic target costs.

#### IV. CONCLUSION

##### SUMMARY OF STATISTICAL RESULTS

The foregoing statistical analysis suggests that some of the advantages usually attributed to incentive contracts may be illusory. It is commonly believed that incentive contracts provide substantial entrepreneurial motivation for increased efficiency and tighter cost control. This belief is one of the stronger justifications for the current extensive use of cost-incentive contracts. The evidence presented here, however, implies that the incentive effect on contractors' costs and efficiency may be weaker than is customarily believed. Rather, the evidence suggests that the cost underruns commonly observed for Air Force incentive contracts are the result of a general upward shift in target costs.

There is an important implication here for improving the effectiveness of incentive contracts. What is needed to make cost-incentive contracts work effectively are tighter target costs. To insure that incentive contracts motivate contractors toward increased efficiency and lower costs, it is essential that the target cost be a realistic estimate of expected actual costs. Thus, future gains in incentive contracting are going to come through improved methods of determining target costs rather than through more elaborate incentive sharing arrangements. Emphasis must be placed on obtaining better target cost information rather than on higher sharing rates and more complex incentive structures.

##### ALTERNATIVES FOR IMPROVING INCENTIVE CONTRACTING

Provided target costs are determined competitively, there is little chance of obtaining targets that significantly exceed contractors' anticipated costs. In the present procurement environment, however, target costs for most of the incentive contracts awarded for major weapon systems are negotiated without benefit of competition (incentive contracts, in fact, often seem to be regarded as a substitute for competition). This is because the DoD typically awards production and follow-on contracts to the original development contractor without competition from alternative suppliers. As a result, effective price rivalry can exist only at the first stage of the program — the development stage. Once the contractor obtains the initial development contract, he is virtually assured of receiving subsequent production and follow-on contracts without fear of competition from other potential producers. Because the targets for these contracts must be negotiated without market price information, it is extremely difficult for the Government to determine whether the resulting target cost is reasonably close to the contractor's



expected cost. Contractors may thus be able to obtain targets sufficiently above their anticipated costs so that the likelihood of achieving a cost underrun and greater profits is increased substantially.\*

One obvious way to determine realistic prices for major weapon systems and also to provide targets that result in real efficiency incentives would be to utilize competition more extensively in weapon system procurements. Of course the extent to which this is possible depends on the nature of the program; that is, on the degree of uncertainty and other program characteristics. Nonetheless, several promising strategies for increasing competition have been proposed in recent years. These techniques range from total package procurement where one contract is awarded competitively for the entire program, to complete separation where each development, production, and follow-on contract is awarded competitively to the same or to various contractors. While one of these alternatives may be more suitable than another for a particular procurement situation, all offer important advantages over present weapon system procurement practices.†

The importance of using competition to determine target costs in weapon system procurements cannot be overemphasized. Nonetheless, there will be many situations in which price rivalry cannot be effectively used — situations where technical uncertainties are large, the number of potential suppliers limited, and so forth. And it does seem likely that a large portion of all weapon system procurements will continue to be made without benefit of competition. In such cases the DoD must rely upon its cost estimating capability to determine reasonable target costs. Thus, another apparent method for increasing the effectiveness of incentive contracts is through improved cost analysis and estimating techniques.

Recognizing the importance of improved cost information, the DoD has given considerable attention to improving its cost estimating capability. They have devoted much effort to developing a comprehensive data base consisting of cost information from previous weapon system acquisitions. The DoD has also improved its cost estimating methodology and its cost reporting system,‡ and some procurement officials now believe that cost estimating techniques can be refined to the point where they become an effective substitute for price competition in establishing realistic target costs.§

---

\*Procurement officials recognize the difficulty accompanying this method of awarding contracts for major weapon systems. For example, in an address before the Institute of Management of Pre-Development Phase of Government Contracts (September 1965) Deputy Assistant Secretary of Defense (Procurement) John M. Malloy stated:

While most production and support contracts are either fixed-price or contain incentives, these arrangements are negotiated for the most part in a noncompetitive environment and may or may not have resulted in the establishment of targets which provide a contractor real and meaningful incentives. These circumstances provide the strongest incentive to increase the competitive aspects of systems procurement.

Nonetheless, none of the more favorable techniques available has been utilized extensively. †Possible techniques include total package procurement, parallel research and development, second sourcing, and separation. G. R. Hall and R. E. Johnson discuss the merits and limitations of these alternatives in "Competition in the Procurement of Military Hard Goods," The RAND Corporation, P-3796 (Mar. 1968).

‡The Truth-in-Negotiations Act (PL 87-653) is intended to insure the reliability and accuracy of contractor-furnished cost information.

§The rationale for this is made clear in the following remarks presented by Harold Asher, former Deputy for Cost Analysis to the Assistant Secretary of Defense (Systems Analysis), in an address to the Operations Research Society of America, October 16, 1966:

"... the assumption is made that DoD is able to estimate the cost of a new weapon system at least as accurately as any single contractor. The reasonableness of this assumption should be apparent. DoD's cost experience is based on all the weapons produced for DoD, while a single company has only its own past programs as an experience base. The assumption is predicated on the effort we are now making to exploit this greater amount of data and experience."



Although cost estimation plays an important role in obtaining improved cost information, it cannot provide cost estimates that are in any sense equivalent to figures that would result through competition among potential suppliers. There are two reasons for this. First, cost estimation relies extensively on past experience to provide estimates of the costs of proposed weapon systems; consequently, such estimates can be no better than the underlying data upon which they are based. If the costs for the previous weapon system procurements were not obtained competitively, the resulting estimates obviously would not be comparable to competitively determined costs. Unfortunately, the majority of weapon system contracts contained in the DoD's data bank were not awarded competitively; in fact many were CPFF, so that costs were possibly several times larger than they might have been otherwise.

Second, even if all contracts included in the data bank had been awarded competitively, the resulting cost estimates would not be equivalent to competitively determined costs. The reason is that cost estimation utilizes data from a number of contracts with different contractors to project the cost of a proposed weapon system. Because some contractors are more efficient than others, this estimated cost is, in reality, an average cost — an estimate of the cost that would result for a firm of average efficiency. As a result, competitively determined costs would generally be lower than estimated costs and the difference could be substantial. Nonetheless, estimated target costs can still provide some positive efficiency incentives for the less efficient contractors and, as a result, are useful in situations where competition is impractical.

In short, although competition is the preferred means for obtaining cost targets, cost estimation provides a useful tool when competition cannot be utilized effectively. The important point is that these estimated costs may be considerably larger than competitively determined costs and might not provide the strongest efficiency incentives. Since competition is probably not feasible in the majority of weapon system procurements, any improvements that can be made in cost-estimating methodology are probably well worthwhile.

Given these constraints, the effectiveness of incentive contracting could presently be improved by utilizing these contracts more selectively. In the past, incentive contracts were applied in numerous cases in which the technical uncertainties were so large that they precluded any meaningful target cost determination. It is important to recognize these situations and either rely on some other form of pricing arrangement or postpone negotiating the target cost until the uncertainty has been resolved. Better project definition prior to negotiating the incentive structure could contribute much toward improving the effectiveness of these contracts.

## SOME FINAL OBSERVATIONS

Nothing can be said here about the total cost of a weapon system under an incentive contract as compared to that under a cost-reimbursable contract. There is no way to analyze how the choice of contract type affects the overall cost of a weapon system; the results obtained here relate only to differences between actual and target costs. The main point demonstrated here is that incentive contracts probably are not saving the Government much money through increased efficiency and better cost control. Consequently, the merits of incentive contracts will have to be judged on other grounds.

Incentive contracts have several important advantages that should not be overlooked. Because of the upward shift in target costs, incentive contracts provide the Government with better program cost information than do cost-reimbursable contracts. Because target costs are more realistic for incentive contracts, they permit better financial planning and budgetary control while eliminating the large overruns characteristic of cost-reimbursable contracts.

Moreover, incentive contracts may have made both the Government and defense contractors a little more cost-conscious than before. Contractors probably have different attitudes toward costs since the advent of incentive contracts than previously, and the Government may be taking the role of a cost-conscious buyer rather than a benevolent sponsor. Consequently, it is possible that these contracts may have resulted in some indirect cost savings. Unfortunately, these salutary effects cannot be measured and quantified.

\* \* \*



# INVENTORY CONTROL OF BY-PRODUCTS\*

Richard V. Evans

*Case Western Reserve University*

*Cleveland, Ohio*

## ABSTRACT

The system to be controlled produces  $n$  products simultaneously in fixed proportions every time it is activated. Demands for the products in any period are components of an  $n$  dimensional vector random variable with known distribution function. Cases of excess demands backlogged and excess demands lost are considered. In the former the notion of  $k$  convexity can be generalized to guarantee relatively simple form for the optimal policy in an  $n$  decision problem. In the latter, this generalization was not successful although when there is no setup cost, a convexity argument can be used to show that the optimal policy has a simple form.

## INTRODUCTION

There is much to be done to breach the gap between inventory theory and management practice. The theory must concern realistic production systems which serve real customers. This can only be done at the expense of introducing more parameters into the description of the system studied. Unfortunately as a consequence the variety of situations grows combinatorially. Each analysis of a new model becomes less distinguishable. Hopefully, as the number of systems which are understood increases, it will become possible to summarize many models in fewer but more abstract ones. The purpose of the present discussion is to suggest that it is possible to extend the study of complex production systems without serious innovation in analysis. At the same time it is hoped that the discussion of difficulties in extending the analysis will stimulate others to study some of these problems.

The model considered here describes one of the possible ways in which several products may be related because of the system which produces them. A situation in which the production system is composed of resources which are limited in amount, but which may be allocated to produce any desired product mix has been studied previously; however, competing for scarce resources is not the only production relationship which occurs (see Ref. [2]). An extreme in another direction is a system in which products are produced simultaneously by the same equipment; the products are by-products of each other. These situations may vary all the way from a rather symmetric case to that of a major product and unimportant by-products. In the rather symmetric case it is reasonable to consider how to produce given production, holding, and penalty costs. As asymmetries become pronounced, it is natural to suggest the possibility of dumping excess inventory of the by-product and thinking of the production decision as though only the major product were involved. Dumping may mean just a price reduction or in more extreme situations physical destruction of the unwanted goods to avoid excessive

---

This work was supported by the National Science Foundation under Grant GK-1333.

holding charges. Even in the relatively symmetric case it is reasonable to examine the physical production process to see if the relative proportion of output can't be modified to help counteract temporary inventory imbalance.

### MODEL

To begin the study of these processes consider the following dynamic programming problem:

$$(1) \quad f_1(x) \equiv 0,$$

$$(2) \quad f_n(x) = \min_{\theta \geq 0} \{C(\theta z) + L(x + \theta z) + \alpha E(f_{n-1}(T(x + \theta z, d)))\},$$

or equivalently

$$(3) \quad f_n(x) = \min_{\substack{y=x+\theta z \\ y \geq x}} \{C(y - x) + L(y) + \alpha E(f_{n-1}(T(y, D)))\},$$

where

$\theta$  is the production level.

$y$  is the inventory after ordering.

$x$  is the inventory before ordering.

$z$  is the contribution to inventory of one unit of production

$$(z > 0); \sum_{i=1}^n = j = 1.$$

$L$  is the expected one period costs as a function of the inventory after ordering. It is assumed convex. For simplicity the demands  $D$  are assumed to be independent identically distributed random variables and thus the period index has been suppressed.

$T$  is the inventory after sales in the period. This function is also assumed the same in all periods.

$\alpha$  is a discount factor ( $\alpha < 1$ ).

$C(\theta z)$  is the production cost and has the form

$$C(\theta z) = \begin{cases} 0 & \text{if } \theta = 0 \\ c\theta + k & \text{if } \theta > 0 \end{cases}.$$

Although  $\theta$  is a single real number,  $x, y, z$  are vectors in a space whose dimension is the number of products. Also vectors are the random variable  $D$  and its values which will be denoted by  $d$ . It is important to note that the situation considered here contrasts with one dimensional inventory problems where it is convenient to use product units to measure production



and inventory levels. In this model the production level,  $\theta$ , might be measured in some natural production unit such as in man hours or machine hours. When one speaks of an inventory level in this model, he means a list of the number of units of each product on hand and this cannot be measured in production units. The conversion of production units into contribution to inventory level is given by the vector  $z$  whose coordinates are the number of units of each product respectively produced by running the production system for our production unit. One may of course convert any feasible addition to inventory say the vector  $y$  to production units by expressing  $y$  as  $\theta z$  some value of  $\theta$  since it is only possible to make additions of this nature under the assumptions of a strict by-product relationship. Finally, when more specific questions about  $L$  arise, it will also be assumed that

$$L(x + \theta z) = E(g(x + \theta z - D)).$$

In all cases the function  $g$  is such that  $dg(x + \theta z - d)/d\theta < -c$  for all points such that all coordinates of  $(x + \theta z - d)$  are negative. This assumption guarantees that it is desirable to engage in this business at least when sale of all products is guaranteed and  $k$  is zero. A stronger assumption would be to guarantee that a unit of production is desirable even when only one product will be sold and the remaining component products must be added to inventory. In many ways the intermediate region of simultaneous positive and negative levels of inventory, which are of course not present in one product theory, gives one of the major qualitative changes which occurs as one shifts from single product to multiproduct theory. The present model with its single production variable but vector random variable for demand is neither fish nor fowl. It is interesting to consider both since this type of production is not uncommon and further the analysis may help in the transition towards studying more complex systems. It is not surprising that there is a simple case in which the analysis concerns a family of one dimensional processes rather than a truly multidimensional one.

The most obvious conclusion which one would like to substantiate is that the optimal policy in period  $n$  has the form

$$(4) \quad \theta_n^*(x) z = \begin{matrix} y_n^*(x) - x & x \leq y_n'(x) \\ 0 & x \geq y_n'(x) \end{matrix},$$

where  $y_n^*(x) = y_n^*(x + Bz)$  and  $y_n'(x) = y_n'(x + Bz)$  for all real numbers  $B$ . Clearly this is the simplest possible policy which one can consider for this problem and the natural generalization of the  $(s, S)$  type policy to this situation. Along any line in the  $z$  direction through the inventory level  $x$  there are two critical numbers  $y_n'$  and  $y_n^*$  such that if  $x \leq y_n'$  then one produces so that the inventory level is raised to  $y_n^*$ . Unfortunately a single proof of this cannot be given for all interesting  $T$  or at least this author has been unable to do so. Thus several cases will be considered.

## BACKLOGGING

For the case of backlogging where  $T(y, d) = y - d$ , the proof that the optimal policy has the form specified in (4) is straight forward. One need only introduce the notion of a function  $f$

being  $k$  convex in the  $z$  direction, meaning that  $f(\delta z)$  is a one dimensional  $k$  convex function  $[3]^*$  of the scalar variable  $\delta$  even though  $f$  is defined over a higher dimensional space of which  $z$  is an element. The use of this concept proof is almost immediate after one establishes the obvious properties of functions of this class. A function which is the average of functions  $k$  convex in the  $z$  direction also has this property. The sum of a convex function and a function which is  $k$  convex in the  $z$  direction is  $k$  convex in the  $z$  direction. Proofs of these properties may be obtained from their one dimensional counterparts with no difficulty.

Although the establishment of the main result involves no tricks, it does require the exploration of some interesting aspects of this system. As usual in studying the  $n$  period problem, one begins with the simple situation of the last period. The convexity of  $L$  and  $C(\theta z)$  for  $\theta > 0$  in the  $z$  direction imply that there is a unique minimum cost in the  $z$  direction which is achieved at a point say  $y_1^*(x)$  for any starting point  $x$ . Moreover, it further guarantees a unique point  $y_1'(x)$  such that  $k + c\theta^*(y_1'(x)) + L(y_1^*(x)) = L(y_1'(x))$ . Thus the optimal policy for the last period has the prescribed form. Because of the assumption that the marginal worth of an extra unit of production is at least as great as its cost when all products are at negative levels and the fact that  $z > 0$ , one can say that for no  $x$  is  $y_1^*(x) \leq 0$ . More than this can only be said if stronger assumptions on  $L$ , or equivalently on  $g$ , are made. The form of the one period optimal expected costs is

$$f_1(x) = \begin{cases} C(y_1^*(x) - x) + L(y_1^*(x)) & \text{for } x < y_1'(x) \\ L(x) & \text{for } x > y_1'(x) \end{cases}.$$

This function is  $k$  convex in the  $z$  direction since this is true for the two pieces separately and no difficulties arise in moving from one piece to the next. The latter is true since for  $x \leq y_1'(x)$   $f_1$  is linear with derivative  $-c$  in the  $z$  direction and  $y_1^*(x) = x + \theta^*(x)$  minimizes  $c\theta(x) + L(x + \theta(x))$ , while  $L(y_1'(x)) = k + c\theta^*(y_1'(x)) + L(y_1^*(x) + \theta^*(y_1'(x))z) = c(y_1^*(x) - y_1'(x)) + L(y_1^*(x))$ . Now in the special case of backlogging in which  $T(y, d) = y - d$  for all demands  $d$ ,  $E(f_1(T(y, d)))$  is  $k$  convex in the  $z$  direction and this is not disturbed by adding the convex function  $L(y)$ . Thus the optimal two period policy has the same form depending on two critical parameters for each line in the  $z$  direction, and an induction may continue the analysis.

#### EXCESS DEMANDS LOST

When a more general  $T$  is used, trouble starts. Possible difficulties are easily identified in developing the analysis of what intuitively should be the simplest of these situations. Suppose that there is no setup cost. Moreover, consider only two products with linear holding and penalty costs with parameters  $h_1, h_2, p_1, p_2$ , respectively. This guarantees not only the convexity of  $g$ , but furthermore that its mixed partial derivatives are zero. For  $T$  assume that there is no backlogging so that  $T(q) = (\max(q_1, 0), \max(q_2, 0))$ . There is of course no difficulty in the last period in establishing the single critical number nature of the optimal policy. The trouble starts, however, immediately since  $E(f_1(T(y, D)))$  must have desirable properties if the policy is to remain simple.

\* $f(x)$  is  $k$  convex if  $f(x) + af'(x) - k - f(x + a)$ .

The situation in the last period is rather special and obvious once one examines the general step in the induction. From the assumption that the  $n-1$  st period policy is of the simple form the optimal discounted costs are

$$f_{n-1}(x) = \begin{cases} c\theta^* + L(\theta^*z + x) + \alpha E(f_{n-2}(T(\theta^*z + x, D))) & \text{for } x < \theta^*(x) \\ L(x) + \alpha E(f_{n-2}(T(x, D))) & \text{for } x > \theta^*(x) \end{cases},$$

where  $\theta^*(x)z + x$  is constant for all  $x$  on a given line in the  $z$  direction. One must now consider  $\alpha E(f_{n-1}(T(y, D)))$  which for simplicity will be denoted by  $E(\Gamma(y - D))$ . Since  $\Gamma$  is constant for  $y - d$  in the negative quadrant, it is obvious that  $\Gamma$  is not convex. Thus one does not have the convexity of the  $n$ th period objective function  $c\theta + L(\theta z + x) + E(\Gamma(\theta z + x - D))$  from the average of convex functions is convex theorem. To use that theorem one must combine terms and consider  $L(y) + \alpha E(\Gamma(y - d))$  as  $E(g(y - D) + \Gamma(y - D)) = E(G(y - D))$  letting  $G = g + \Gamma$ . For further notational convenience let  $E(G(y - D)) = J(y)$ .

For  $y - d$  in the negative quadrant  $G$  is convex because  $g$  is convex and  $\Gamma$  is constant. The first partial derivatives are  $-p_1$  and  $-p_2$ , respectively, while the mixed partial derivative is zero as are the second partials. The positive quadrant is also straight forward inheriting properties directly from  $f_{n-1}$ . Again  $g$  is linear and its partial derivatives are now  $h_1$  and  $h_2$ , respectively. Convexity follows as long as the induction can maintain the convexity of  $f_{n-1}$  with partials not less than  $-p_1$  and  $-p_2$ , respectively, while  $z_1 f_{n-1} + z_2 f_{n-2}$  must be not less than  $-c$ . It is also important that the mixed partial be nonpositive. In the two remaining quadrants the analysis is of course symmetric so for definiteness consider the situation in which  $y_1 - d_1 \geq 0$  and  $y_2 - d_2 \leq 0$ . Here  $g$  is linear as always with partials of  $h_1$  and  $-p_2$ , respectively. In this quadrant  $T(y, d) = (y_1 - d_1, 0)$  and so  $\Gamma$  is constant with respect to the second coordinate so there is no problem with convexity in this direction nor the mixed partial which must be zero. It is the behavior relative to the first coordinate which is interesting, although with these assumptions one has convexity here from  $f_{n-1}$ . The partial derivatives also give no trouble so that in putting the four pieces together one gets a function convex everywhere. The expectation is then convex and thus so is the function to be minimized by the choice of the policy for period  $n$ .

Now the single critical number aspect of the policy is an immediate consequence of the convexity. Thus

$$f_n(x) = \begin{cases} c\theta^*(x) + J(\theta^*(x)z + x) & 0 \leq \theta^*(x) \\ J(x) & \theta^*(x) \leq 0 \end{cases}.$$

Convexity is immediate in the region in which no ordering is done and continuity at the boundary is also obvious. In the ordering region

$$\begin{aligned} f_{n-1} &= c\theta_1^* + J_1(\theta_1^*z + x) \left[ \theta_1^*z_1 + 1 \right] + J_2(\theta_1^*z + x) \left[ \theta_1^*z_2 \right] \\ &= \left[ c + J_1z_1 + J_2z_2 \right] \theta_1^* + J_1. \end{aligned}$$

The optimality condition is that the bracketed term in this expression be zero thus

$$f_{n-1} = J_1 (\theta^* z + x)$$

and

$$z_1 f_{n-1} + z_2 f_{n-2} = -c.$$

Differentiating this expression with respect to  $x_1$  gives

$$z_1 f_{n-11} + z_2 f_{n-12} = 0$$

$$z_1 f_{n-11} = -z_2 f_{n-12}.$$

Explicitly these partials are:

$$f_{n-11} = J_{11} [\theta_1^* z_1 + 1] + J_{12} [\theta_1^* z_2];$$

$$\begin{aligned} f_{n-12} &= J_{11} [\theta_2^* z_1] + J_{12} [\theta_2^* z_2 + 1] \\ &= J_{12} [\theta_1^* z_1 + 1] + J_{22} [\theta_1^* z_2]; \end{aligned}$$

$$f_{n-22} = J_{12} [\theta_2^* z_1] + J_{22} [\theta_2^* z_2 + 1].$$

For convexity first one requires non-negativity of the second partials and thus one must examine the partials of  $\theta^*(x)$ . This can be done by differentiating the optimality condition giving

$$J_{11} z_1 [\theta_1^* z_1 + 1] + J_{12} z_1 [\theta_1^* z_2] + J_{21} z_2 [\theta_1^* z_1 + 1] + J_{22} z_2 [\theta_1^* z_2] = 0$$

or

$$\theta_1^* = \frac{-[J_{11} z_1 + J_{12} z_2]}{[J_{11} z_1^2 + J_{12} z_1 z_2] + [J_{22} z_2^2 + J_{12} z_1 z_2]}$$

and similarly

$$\theta_2^* = \frac{-[J_{12} z_1 + J_{22} z_2]}{[J_{11} z_1^2 + J_{12} z_1 z_2] + [J_{22} z_2^2 + J_{12} z_1 z_2]}$$

By the induction hypothesis all the brackets are non-negative and thus

$$z_1 \theta_1^* + z_2 \theta_2^* = -1 \quad \text{and} \quad 0 \geq z_1 \theta_1^* \geq -1, \quad 0 \geq z_2 \theta_2^* \geq -1.$$



From this the second partials of  $f_n$  are non-negative

$$\begin{aligned} f_{n-11} &= \left[ z_1 J_{11} + z_2 J_{12} \right] \theta_1^* + J_{11} \\ &\geq J_{11} - \left[ z_1 J_{11} + z_2 J_{12} \right] \\ &= z_2 (J_{11} - J_{12}) \geq 0 \end{aligned}$$

$$\begin{array}{cccc} f_{n-11} & f_{n-12} & z_1 f_{n-11} & z_1 f_{n-12} \\ & & = \frac{1}{z_2 z_1} & \geq 0 \\ f_{n-12} & f_{n-22} & z_2 f_{n-12} & z_2 f_{n-22} \end{array}$$

since

$$z_1 f_{n-11} = -z_2 f_{n-12}$$

and

$$z_2 f_{n-22} = -z_1 f_{n-12}.$$

Clearly there is no difficulty in putting the pieces of  $f$  together because at the boundary of the order region  $f_{n-1} = J_1(x)$  which is also the initial value in the no ordering region, and the negativity of  $\theta_1^*$  and  $\theta_2^*$  imply that the boundary is monotone decreasing in either variable so that there is no possibility that increasing inventory in one product can cause a reentry into the ordering region. Moreover, clearly inductively the partials of  $f_n$  satisfy the inequalities required in the discussion of  $g + \Gamma$ . All that has been said here certainly applies to the analysis of  $f_1$  given the simple nature of the one period losses and an assumption that  $f_0 = 0$ .

## CONCLUSIONS

From these models one may conclude that in several situations by-product type production systems will have optimal policies which are of simple form. In the backlogging case, the argument rests on a straightforward generalization of the notion of  $k$  convexity to  $k$  convexity in the  $z$  direction. In the no backlogging case difficulties arise in that increasing stock in the  $z$  direction does not increase period ending inventories in this direction. Thus when one tries to piece together the expected  $n$  period costs if there is  $y$  on hand after ordering, he finds three pieces. The first is convex and decreases with derivative  $-p_1 z_1 - p_2 z_2$  in the  $z$  direction. The second piece is the bad actor for unless one has a convexity assumption, which is possible when the setup cost is zero, this piece need not be convex although the derivative of  $g + \Gamma$  as  $y$  increases in the  $z$  direction is sufficiently small throughout the region. It appears that one needs a convex function with derivative less than the derivative in the next region in order to combine pieces with a final  $k$  convex piece to get a  $k$  convex function. This problem is intensified in the setup cost no backlogging case because as  $T(y - d)$  moves along an axis it may pass from the ordering region to a region in which on a marginal cost basis it is desirable



to order, but no ordering is done because of the setup cost. This is not possible in the case of no setup cost because here the marginal properties determine the optimal policy completely. The no backlogging setup cost problem can undoubtedly be rescued by a poly frequency function argument with appropriate restrictions on the cost parameters as it can be in the one dimensional situation [1].

Qualitatively there are two phenomenon which are characteristic of multiproduct problems which appear in these situations. The first is that the one dimensional point of no inventory becomes the region of partial stock outs. Similarly, the second is that as one increases the stock after ordering the period ending level  $y - d$  may pass from large negative values to large positive ones without passing through an ordering region in the no backlogging situation. As far as the notion of  $k$  convexity in a given direction is concerned, this is easily shown to be preserved under averaging; however, it is not easy to show that the minimization preserves  $k$  convexity along any path other than one actively involved in the minimization. This causes trouble with the no backlogging case and with the obvious generalization to permitting production moves in any of several specified directions. It is reasonable to ask if these troubles are real or merely specters of ones imagination. So far even in multiple direction problems there are no examples showing that multiple disconnected ordering regions occur. In fact, the few problems the author has solved indicate that if such examples do, in fact, exist, they will be somewhat pathological. Thus, it seems reasonable to continue the search for something beyond  $n$  dimensional convexity which can be maintained inductively and is not disturbed by the minimization operation, averaging, and the transformation from initial to ending inventory.

#### REFERENCES

- [1] Arrow, K. J., S. Karlin, and H. Scarff, Studies in the Mathematical Theory of Inventory and Production (Stanford University Press, Stanford, Calif., 1959).
- [2] Evans, R. V., "Inventory Control of a Multiproduct System with a Limited Production Resource," NRLQ, 14, 173-184 ( ).
- [3] Scarff, H., "The Optimality of  $(S, s)$  Policies in the Dynamic Inventory Problem," Chapter 13 in Arrow, K., S. Karlin, and P. Sappes, (Editors) Mathematical Methods in the Social Sciences 1959 (Stanford University Press, Stanford, Calif., 1960).

\* \* ' \*

# DISCOUNTED PRODUCTION SCHEDULING AND EMPLOYMENT SMOOTHING\*

Steven A. Lippman† and John S. C. Yuan‡

*Stanford University*

## ABSTRACT

We consider the problem of minimizing the sum of production, employment smoothing, and inventory costs over a finite number of time periods where demands are known. The fundamental difference between our model and that treated in [1] is that here we permit the smoothing cost to be nonstationary, thereby admitting a model with discounting. We show that the values of the instrumental variables are nondecreasing in time when demands are nondecreasing. We also derive some asymptotic properties of optimal policies.

## 1. INTRODUCTION AND SUMMARY

We consider the problem of minimizing the sum of production, employment smoothing, and inventory costs over a finite number of time periods. We require that the known demands be met each period. There are two modes of production, termed regular-time and overtime; the latter is constrained to be no greater than a fixed proportion of the work force. We also permit a distinction between the amount of labor employed at regular-time and the level of the work force. There is a piece-wise linear cost associated with fluctuations in the work force, a nondecreasing inventory holding cost, and a production cost which subsumes convexity in the levels of regular-time, overtime, and work force.

The fundamental difference between our model and that which we treated in [1] is that here we permit the smoothing cost to decrease in time, thereby admitting a model with discounting. To our knowledge, this is the only such model in the "production scheduling, employment smoothing" literature. But the added generality creates sizeable analytic difficulties and many of the results in [1] are no longer valid. In Theorem 1 of [1], we showed that the cumulative of regular-time plus overtime production was bounded above by the upper convex envelope of the cumulative demand schedule. In Theorem 3 of [1], we showed that the values of all the instrumental variables were secularly decreasing under the assumptions of zero initial inventory and secularly decreasing demand. A three-period example, however, demonstrates that neither of these results holds for the model in this paper. Assuming that the known demands are nondecreasing in time, we prove in this paper that most of Theorems 2 and 4 of [1] is valid, but we need to give an entirely new proof.

---

\*This work was supported in part by a National Science Foundation Grant (GS-552) and by an Office of Naval Research Contract Nonr 225(77). The authors are grateful to Professor Harvey M. Wagner for helpful suggestions on this paper.

†Now at the Graduate School of Business, University of California, Los Angeles.

‡Now at Management Science Services, IBM Corporation, Armonk, New York.

Unfortunately, even without the possibility of overtime or of idle-time, the main algorithm presented in [2] cannot be salvaged, because it depends critically on stationarity of the smoothing cost.

## 2. MODEL DESCRIPTION

The following is a description of the model and its underlying assumptions. The reader is referred to pp. 3-12 of the previous paper [1] for a more detailed discussion on the meaning and economic implications of the assumptions.

We desire a minimum cost production plan over a finite number of time periods  $T$ . In each period  $t$ ,  $t=1, 2, \dots, T$ , the instrumental variables, measured in hours, are

$$\begin{aligned} u_t &= \text{work force producing at regular time,} \\ z_t &= \text{work force producing on overtime, and} \\ w_t &= \text{total work force.} \end{aligned}$$

For each  $t$  we impose the constraints\*

$$(2.1) \quad 0 \leq u_t \leq w_t$$

and

$$(2.2) \quad 0 \leq z_t \leq \alpha w_t \quad (0 \leq \alpha)$$

so that  $(w_t - u_t)$  represents the idle portion of the work force, and overtime production does not exceed a fixed proportion of the work force.

The fluctuations in the regular-time work force and inventory are implicit decision variables. The hiring and firing smoothing costs associated with period  $t$ ,  $t=1, 2, \dots, T+1$ , are

$$(2.3a) \quad s_t(w_t, w_{t-1}) = \begin{cases} g_t \cdot (w_t - w_{t-1}) & \text{for } w_t > w_{t-1} \\ f_t \cdot (w_{t-1} - w_t) & \text{for } w_t \leq w_{t-1} \end{cases},$$

where  $w_0 = 0$  is the initial work force. We assume

$$(2.3b) \quad f_t, g_t \geq 0 \quad \text{and} \quad f_t \geq f_{t+1}, g_t \geq g_{t+1}.$$

Note that (2.3a) and (2.3b) constitute the only distinction between the model and [1]. Our analysis encompasses two terminal assumptions. If at the end of period  $T$ , the work force,  $w_T$ , is to be fired at the cost  $f_{T+1}$ , then we let  $w_{T+1} = 0$  at  $t = T+1$  in (2.3a). If this work force is not to be fired, or is to be fired at no cost, then we let  $w_{T+1} = w_T$  in (2.3a).

Inventory at the end of period  $t$  is described by

\*The equations of this section are numbered as in [1] for the readers convenience.

$$(2.4) \quad i_t = i_{t-1} + u_t + z_t - d_t = i_0 + \sum_{j=1}^t (u_j + z_j - d_j),$$

where  $d_t \geq 0$  is the associated period demand requirement measured in hours, and  $i_0 \geq 0$  the inventory entering period 1. We assume all demand must be met and backlogging is not permitted, so that we constrain

$$(2.5) \quad i_t \geq 0.$$

We let

$h_t(i)$  = holding cost of ending inventory  $i$  for period  $t$

and assume only that

$$(2.6) \quad h_t(i) \text{ is a nondecreasing function.}$$

Turning to the remaining component of total cost, we let

$\pi_t(u, z, w)$  = total production cost for period  $t$ ,

where  $u_t = u$ ,  $z_t = z$ , and  $w_t = w$  and satisfy (2.1) and (2.2). We make the following assumptions about production costs\*:

$$(2.11) \quad \pi_t(u + \epsilon, z + \delta, w + \beta) - \pi_t(u, z, w) \geq 0 \quad \text{for } \epsilon, \delta, \beta \geq 0,$$

$$(2.12) \quad \pi_t(u + \epsilon, 0, w + \beta) - \pi_t(u, 0, w + \beta) \geq \pi_t(u + \epsilon, 0, w) - \pi_t(u, 0, w) \quad \text{for } \epsilon, \beta \geq 0,$$

$$(2.13) \quad \pi_t(w + \epsilon + \beta, 0, w + \epsilon + \beta) - \pi_t(w + \beta, 0, w + \beta) \geq \pi_t(w + \epsilon, 0, w + \epsilon) - \pi_t(w, 0, w) \quad \text{for } \epsilon, \beta \geq 0,$$

$$(2.14) \quad \pi_t(w + \beta, z, w + \beta) - \pi_t(w, z, w) \geq \pi_t(w + \beta, 0, w + \beta) - \pi_t(w, 0, w) \quad \text{for } \beta \geq 0,$$

$$(2.15) \quad \pi_t(u + \epsilon + \beta, 0, w) - \pi_t(u + \beta, 0, w) \geq \pi_t(u + \epsilon, 0, w) - \pi_t(u, 0, w) \quad \text{for } \epsilon, \beta \geq 0,$$

$$(2.16) \quad \pi_t(w, z + \eta + \delta, w) - \pi_t(w, z + \eta, w) \geq \pi_t(w, z + \delta, w) - \pi_t(w, z, w) \quad \text{for } \eta, \delta \geq 0,$$

$$(2.17) \quad \pi_t(u, z, w) \geq \pi_t(u + \epsilon, z - \epsilon, w) \quad \text{for } \epsilon \geq 0,$$

$$(2.18) \quad \pi_t(w, z, w) \geq \pi_t(w + \epsilon, z - \epsilon, w + \epsilon) \quad \text{for } \epsilon \geq 0,$$

$$(2.19) \quad \pi_t(u + \epsilon, z + \delta, w + \beta) - \pi_t(u, z, w) \geq \pi_{t+1}(u + \epsilon, z + \delta, w + \beta) - \pi_{t+1}(u, z, w) \quad \text{for } \epsilon, \delta, \beta \geq 0.$$

\*For brevity, we do not display the various restrictions on the values of  $\epsilon, \delta, \eta$  in (2.11) through (2.19) and (A1) implied by (2.1) and (2.2), but we assume these restrictions are in force. Further in the proofs below, (2.11) and (2.19) are needed only when  $(z + \delta) [(w + \epsilon) - (u + \epsilon)] = 0$  for  $\epsilon, \delta, \beta \geq 0$ .

We now introduce an additional assumption which is employed only in the proof of part (iv) of Theorem 1; namely,

$$(A1) \quad \pi_t(w, z + \epsilon + \Delta, w) - \pi_t(w, z + \epsilon, w) \geq \pi_t(w + \epsilon, z + \Delta, w + \epsilon) - \pi_t(w + \epsilon, z, w + \epsilon) \quad \text{for } \epsilon, \Delta > 0.$$

To aid in the interpretation of (A1), define  $\pi'_t(z, w) \equiv \pi_t(w, z, w)$  and assume  $\pi'_t(\cdot, \cdot)$  is smooth enough, so that we can rewrite (A1) as

$$\partial \pi'_t(z + \epsilon, w) / \partial z \geq \partial \pi'_t(z, w + \epsilon) / \partial z, \quad \epsilon \geq 0.$$

Then (A1) implies that the marginal cost of overtime is more heavily dependent on the level of overtime than on the level of regular-time. The reasonability of (A1) might be justified by considering  $\partial \pi'_t(z + \epsilon, w) / \partial z$  to be the marginal cost of overtime that accrues when  $w$  workers (in hours) provide  $z + \epsilon$  hours of overtime. Since these  $w$  workers will be more "fatigued" by their labor than  $w + \epsilon$  workers will be by theirs (as  $\frac{w + z + \epsilon}{w} > \frac{w + \epsilon + z}{w + \epsilon}$ ), the  $w$  workers will be less efficient, and, hence, more costly in providing additional labor. If  $\pi_t(u, z, w)$  is separable in  $u$ ,  $z$ , and  $w$ , then (A1) is satisfied.

Let  $U$ ,  $Z$ , and  $W$  be  $T$ -dimensional vectors of real numbers. We refer to the  $3T$ -dimensional vector  $(U, Z, W)$  as policy  $(U, Z, W)$  and we say that  $(U, Z, W)$  is feasible if it satisfies (2.1), (2.2), (2.4), and (2.5). The problem, then, is to minimize

$$Q(U, Z, W) \equiv \sum_{t=1}^T \pi_t(u_t, z_t, w_t) + \sum_{t=1}^{T+1} s_t(w_t, w_{t-1}) + \sum_{t=1}^T h_t(i_t)$$

subject to (2.1), (2.2), (2.4), (2.5), and where the cost functions satisfy (2.3a), (2.3b), (2.6), (2.11) through (2.19).

### 3. FORM OF OPTIMAL POLICIES UNDER INCREASING DEMANDS

We now characterize the form of an optimal policy under the assumption of increasing demands.

**THEOREM 1:** If

$$(3.1) \quad w_0 \leq \max(d_1 - i_0, 0)$$

$$(3.2) \quad d_t \leq d_{t+1} \quad \text{for } t = 1, 2, \dots, T - 1,$$

then there exists an optimal policy  $(U, Z, W)$  such that

$$(i) \quad w_t \leq w_{t+1} \quad \text{for } t = 0, 1, 2, \dots, T - 1$$

$$(ii) \quad u_t = w_t \quad \text{for } t = 1, 2, \dots, T$$

$$(iii) \quad u_t \leq u_{t+1} \quad \text{for } t = 1, 2, \dots, T - 1.$$



Moreover, if (A1) holds, then  $(U, Z, W)$  also satisfies

$$(iv) \quad u_t + z_t \leq u_{t+1} + z_{t+1} \leq d_T \quad \text{for } t=1, 2, \dots, T-1.$$

If assumption (3.1) is dropped, then  $w_t$  may decrease, as is shown by the following example:

$$d_1 = d_2 = 1, \quad T = 2, \quad w_0 = 1 + \delta, \quad i_0 = 0, \quad h_1(i) = i,$$

$$f_1 = 2, \quad f_2 = 1 - \epsilon, \quad \alpha = 1, \quad \delta > 0, \quad 0 < \epsilon < 1,$$

$$\pi_t(u, z, w) = w + 2z \quad \text{for } t=1, 2.$$

The unique optimal policy, regardless of whether  $w_{T+1} \equiv w_T$  or  $w_{T+1} = 0$ , is  $u_1 = u_2 = 1$ ,  $w_1 = 1 + \delta$ ,  $w_2 = 1$ ,  $z_1 = z_2 = 0$ ; i.e.,  $w_1 > w_2$ .

In [1] we showed the existence of an optimal  $(U, Z, W)$  that satisfied  $z_t \leq z_{t+1}$ ,  $t=1, 2, \dots, T-1$  in addition to (i)-(iv). With nonstationary smoothing costs, however, this is not the case, as can be seen by the following example:

$$d_1 = 2, \quad d_2 = 4, \quad T = 2, \quad f_1 = f_2 = 0,$$

$$g_1 = 4, \quad g_2 = 1, \quad i_0 = 0, \quad h_1(i) = h_2(i) = 0,$$

$$\alpha = 1, \quad \pi_t(u, z, w) = 3z + 1w \quad \text{for } t=1, 2.$$

The unique optimal policy is  $u_1 = w_1 = 1$ ,  $z_1 = 1$ ,  $u_2 = w_2 = 4$ ,  $z_2 = 0$ ; so that  $z_1 \geq z_2$ .

#### 4. PROOF OF THEOREM 1

We have subdivided the proof of (i), which entails numerous details, in order to highlight the main arguments and constructions. We start by proving six preliminary lemmas which will be employed in the proof of (i). The reader should note that the construction given in the proof of (i) is not finite, and that the nonstationary smoothing cost renders the proof of Theorem 2 in [1] unusable, thereby necessitating this entirely new proof.

We use the following lemma for comparing smoothing costs.

LEMMA 1: Let  $w_t \geq 0$ ,  $t = 0, 1, 2, \dots, T$ , have values such that for  $0 \leq j \leq k \leq T-1$ ,

$$w_j \leq w_{j+1} = \dots = w_{k-1} = w_k \quad \text{and} \quad w_k \geq w_{k+1}.$$

Further, let  $w'_t$  be such that

$$w'_t = w_t \quad \text{for all } t \neq j+1, \dots, k+1, \quad w_j \leq w'_{j+1} = \dots = w'_k \leq w_k, \quad \text{and} \quad w'_k \geq w'_{k+1} \geq w_{k+1}.$$

Then

$$S(W) \geq S(W') \quad \text{and} \quad S^+(W) \geq S^+(W')$$

where

$$S(W) = \sum_{t=1}^T s_t(w_t, w_{t-1}) \quad \text{and} \quad S^+(W) = \sum_{t=1}^{T+1} s_t(w_t, w_{t-1}).$$

The proof of Lemma 1 is a straightforward application of (2.3b).

Let  $(\tilde{U}, \tilde{Z}, \tilde{W})$  be any policy. If  $\tilde{w}_t \leq \tilde{w}_{t+1}$  for  $t=0, 1, 2, \dots, T-1$ , let  $\tilde{r} = T$ . Otherwise, let

$$\tilde{r} = \min \{t: \tilde{w}_t > \tilde{w}_{t+1}, 0 \leq t \leq T-1\}.$$

We now define  $S$  = set of policies which satisfy the feasibility constraints (2.1), (2.2), (2.4), (2.5).

$$S(1, r) = \{(\tilde{U}, \tilde{Z}, \tilde{W}) \in S: \tilde{r} \geq r\} \quad \text{for } r=1, 2, \dots, T.$$

$$S(2, r) = \{(\tilde{U}, \tilde{Z}, \tilde{W}) \in S(1, r): \tilde{z}_{r+1} = 0\} \quad \text{for } r=1, 2, \dots, T-1.$$

$$S(3, r) = \bigcup_{m=0}^{r-1} \{(\tilde{U}, \tilde{Z}, \tilde{W}) \in S(2, r): \tilde{w}_m < \tilde{w}_{m+1} = \dots = \tilde{w}_r\} \\ \text{for } r=1, 2, \dots, T-1.$$

$$S(4, r) = \bigcup_{m=0}^{r-1} \{(\tilde{U}, \tilde{Z}, \tilde{W}) \in S(3, r): \tilde{u}_{m+1} = \dots = \tilde{u}_r = \tilde{w}_r\} \\ \text{for } r=1, 2, \dots, T-1.$$

$$S(5, r, m) = \{(\tilde{U}, \tilde{Z}, \tilde{W}) \in S(4, r): \tilde{z}_t = 0, m+1 \leq t \leq r+1\} \\ \text{for } m=0, 1, \dots, r-1 \text{ and } r=1, 2, \dots, T-1.$$

$$S(5, r) = \bigcup_{m=0}^{r-1} S(5, r, m) \quad \text{for } r=1, 2, \dots, T-1.$$

The relation between Lemmas 2 through 6 and the sets  $S(i, r)$  is summarized in Figure

1. A branch, denoted by  $A \xrightarrow{Li} B$ , means that given a policy in set A, Lemma i ensures that



there be a less-cost policy in either set B or set C. For fixed  $r$  we have

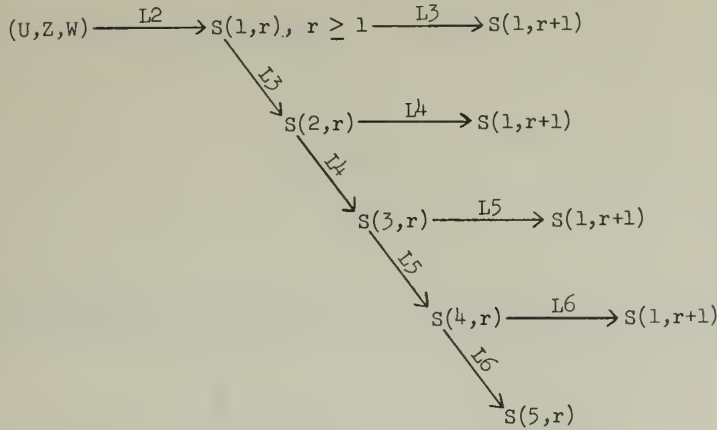


Figure 1. Diagrammatic representation of Lemmas 2-6

$$S \supset S(1, r) \supset \dots \supset S(5, r).$$

Hence, Lemmas 2 through 6 show that starting from any feasible policy and its corresponding  $r$ , we can find a less-cost policy that either is in  $S(5, r)$  or in  $S(1, r+1)$ .

**LEMMA 2:** There is a policy  $(U', Z', W') \in S(1, 1)$  such that  $Q(U', Z', W') \leq Q(U, Z, W)$  for all  $(U, Z, W) \in S \sim S(1, 1)$ .

**PROOF:** Take  $(U, Z, W) \in S \sim S(1, 1)$ , so that  $w_0 > w_1$ . Let  $(U', Z', W') \in S(1, 1)$  be defined by

$$u'_t = u_t, w'_t = w_t, z'_t = z_t \quad \text{for } t=2, 3, \dots, T$$

and

$$u'_1 = w'_1 = w_0, z'_1 = u_1 + z_1 - u'_1.$$

As  $i_t = i'_t$  for all  $t$ , the inventory cost of the two policies is the same. Letting  $j=k=0$ , we have that  $w_t$  and  $w'_t$  satisfy the hypotheses of Lemma 1, and hence both  $S(W') \leq S(W)$  and  $S^+(W') \leq S^+(W)$ .

Letting  $\epsilon = w_0 - u_1 > 0$ , we have from (2.11) and (2.18) that

$$\pi_1(u_1, z_1, w_1) \geq \pi_1(u_1, z_1, u_1) \geq \pi_1(u_1 + \epsilon, z_1 - \epsilon, u_1 + \epsilon) = \pi_1(u'_1, z'_1, w'_1).$$

Thus  $Q(U', Z', W') \leq Q(U, Z, W)$ .

Q.E.D.

**LEMMA 3:** There is a policy  $(U', Z', W') \in S(2, r) \cup S(1, r+1)$  such that  $Q(U', Z', W') \leq Q(U, Z, W)$  for all  $(U, Z, W) \in S(1, r) \sim [S(2, r) \cup S(1, r+1)]$ ,  $r=1, 2, \dots, T-1$ .

**PROOF:** Take  $(U, Z, W) \in S(1, r) \sim [S(2, r) \cup S(1, r+1)]$  so that  $w_r > w_{r+1}$  and  $z_{r+1} > 0$ . Let  $(U', Z', W')$  be defined by

$$u'_t = u_t, w'_t = w_t, z'_t = z_t \quad \text{for } t \neq r+1;$$

and

$$u'_{r+1} = u_{r+1} + z_{r+1}, w'_{r+1} = \max\{u'_{r+1}, w_{r+1}\}, z'_{r+1} = 0 \quad \text{if } u_{r+1} + z_{r+1} \leq w_r,$$

$$u'_{r+1} = u_{r+1} + (w_r - w_{r+1}), w'_{r+1} = w_r, z'_{r+1} = z_{r+1} - (w_r - w_{r+1}) \quad \text{if } u_{r+1} + z_{r+1} > w_r.$$

In either case, the inventory costs are equal while  $S^+(W') \leq S^+(W)$  and  $S(W') \leq S(W)$  by Lemma 1 with  $j = k = r$ . Also  $\pi_{r+1}(u_{r+1}, z_{r+1}, w_{r+1}) \geq \pi_{r+1}(u'_{r+1}, z'_{r+1}, w'_{r+1})$  by (2.17) and (2.18), so  $Q(U', Z', W') \leq Q(U, Z, W)$ .

If  $u_{r+1} + z_{r+1} \leq w_r$ , then  $z'_{r+1} = 0$  and so  $(U', Z', W') \in S(2, r)$ .

If  $u_{r+1} + z_{r+1} > w_r$ , then  $r' \geq r+1$  so that  $(U', Z', W') \in S(1, r+1)$ . Q.E.D.

LEMMA 4: There is a policy  $(U', Z', W') \in S(1, r+1)$  such that  $Q(U', Z', W') \leq Q(U, Z, W)$  for all  $(U, Z, W) \in S(2, r) \sim [S(3, r) \cup S(1, r+1)]$ ,  $r = 1, 2, \dots, T-1$ .

PROOF: Take  $(U, Z, W) \in S(2, r) \sim [S(3, r) \cup S(1, r+1)]$ , then

$$(4.1) \quad w_0 = w_1 = \dots = w_r > w_{r+1} \quad \text{and} \quad z_{r+1} = 0.$$

Let

$$w'_t = w_t \quad \text{for } t \neq r+1, w'_{r+1} = w_r;$$

$$u'_t = \min\{u_t + z_t, w_t\} = u_t + \gamma_t \quad \text{for } t \leq r, u'_{r+1} = w'_{r+1}, u'_t = u_t \quad \text{for } t > r+1;$$

$$z'_t = \max \left\{ z_t - \left[ w_r - u_{r+1} - (w_t - u_t) - \sum_{i=t+1}^r (z_i - (z'_i + \gamma_i)) \right], 0 \right\} \quad \text{for } t \leq r;$$

$z'_t = z_t$  for  $t > r$  where  $\gamma_t = \min\{w_t - u_t, z_t\}$  for  $t \leq r$ . The feasibility of  $(U', Z', W')$  follows from  $z'_{r+1} = 0$ , (3.2), and  $i_r \geq w_r - u_{r+1}$ . To show that

$$\sum_{t=1}^{r+1} [\pi_t(u_t, z_t, w_t) - \pi_t(u'_t, z'_t, w'_t)] \geq 0,$$

it suffices to show that

$$\begin{aligned} \pi_t(u_t, z_t, w_t) - \pi_t(u_t + \gamma_t, z_t - \Delta_t - \gamma_t, w_t) &\geq \pi_{r+1} \left( u_{r+1} + \sum_{j=t}^r \Delta_j, 0, w_{r+1} + \sum_{j=t}^r \epsilon_j \right) \\ &- \pi_{r+1} \left( u_{r+1} + \sum_{j=t+1}^r \Delta_j, 0, w_{r+1} + \sum_{j=t+1}^r \epsilon_j \right) \quad \text{where } \gamma_t = \min\{w_t - u_t, z_t\}, \end{aligned}$$

$$\Delta_t = \max\{z_t - z'_t - \gamma_t, 0\},$$

and

$$\epsilon_t = \min \left\{ \Delta_t, w_t - w_{r+1} - \sum_{j=t+1}^r \epsilon_j \right\} \quad \text{for } t=1, 2, \dots, r.$$

Note that  $\gamma_t(\Delta_t)$  is the amount of overtime shifted to regular-time in period  $t$  ( $r+1$ ). The inequality follows immediately from (2.17) if  $\Delta_t = 0$ . If  $\Delta_t > 0$ , then  $u_t + \gamma_t = w_t$  so that

$$\begin{aligned} \pi_t(u_t, z_t, w_t) &= \pi_t(u_t + \gamma_t, z_t - \Delta_t - \gamma_t, w_t) \\ &\geq \pi_t(w_t, z_t - \gamma_t, w_t) - \pi_t(w_t, z_t - \gamma_t - \Delta_t, w_t) && \text{by (2.17)} \\ &\geq \pi_{r+1}(w_t, z_t - \gamma_t, w_t) - \pi_{r+1}(w_t, z_t - \gamma_t - \Delta_t, w_t) && \text{by (2.19)} \\ &\geq \pi_{r+1}(w_t, \Delta_t, w_t) - \pi_{r+1}(w_t, 0, w_t) && \text{by (2.16)} \\ &\geq \pi_{r+1}(w_t + \Delta_t, 0, w_t + \Delta_t) - \pi_{r+1}(w_t, 0, w_t) && \text{by (2.18)} \\ &\geq \pi_{r+1} \left( u_{r+1} + \sum_{j=t}^r \Delta_j, 0, w_t + \Delta_t \right) - \pi_{r+1} \left( u_{r+1} + \sum_{j=t+1}^r \Delta_j, 0, w_t \right) && \text{by (2.15)} \\ &\text{as } w_t = w_r = u_{r+1} + \sum_{j=1}^r \Delta_j \\ &\geq \pi_{r+1} \left( u_{r+1} + \sum_{j=t}^r \Delta_j, 0, w_{r+1} + \sum_{j=t}^r \epsilon_j \right) \\ &\quad - \pi_{r+1} \left( u_{r+1} + \sum_{j=t+1}^r \Delta_j, 0, w_{r+1} + \sum_{j=t+1}^r \epsilon_j \right) && \text{by (2.11)} \\ &\text{as } \Delta_t \geq \epsilon_t \quad \text{and (2.12).} \end{aligned}$$

Thus it follows from the above, Lemma 1 (with  $j = k = r$ ), and (2.6) that  $Q(U', Z', W') \leq Q(U, Z, W)$ . Note that  $(U', Z', W') \in S(1, r+1)$ . Q.E.D.

**LEMMA 5:** There is a policy  $(U', Z', W') \in S(4, r) \cup S(1, r+1)$  such that  $Q(U', Z', W') \leq Q(U, Z, W)$  for all  $(U, Z, W) \in S(3, r) \sim [S(4, r) \cup S(1, r+1)]$ , for  $r=1, 2, \dots, T-1$ .

**PROOF:** Take  $(U, Z, W) \in S(3, r) \sim [S(4, r) \cup S(1, r+1)]$  so that there is an  $m$  and  $p$ ,  $0 \leq m \leq r-1$ ,  $1 \leq p \leq r-m$  such that



$$w_m < w_{m+1} = \dots = w_r = u_{m+1} = \dots = u_{m+p-1} > w_{r+1} \text{ and } \Delta \equiv w_{m+p} - u_{m+p} > 0.$$

By (2.17) we may assume

$$(4.2) \quad z_t > 0 \text{ implies } u_t = w_t \quad \text{for all } t.$$

We may also assume (else a new policy  $(\tilde{U}, \tilde{Z}, \tilde{W})$  could be defined with  $\tilde{p} > p$ )

$$(4.3) \quad z_t = 0 \quad \text{for } m+1 \leq t \leq m+p,$$

as

$$\begin{aligned} & \pi_t(w, z, w) - \pi_t(w, z - \epsilon, w) \\ & \geq \pi_{m+p}(w, z, w) - \pi_{m+p}(w, z - \epsilon, w) \\ & \geq \pi_{m+p}(w, \epsilon, w) - \pi_{m+p}(w, 0, w) \\ & \geq \pi_{m+p}(w + \epsilon, 0, w + \epsilon) - \pi_{m+p}(w, 0, w + \epsilon) \\ & \geq \pi_{m+p}(u + \epsilon, 0, w + \epsilon) - \pi_{m+p}(u, 0, w + \epsilon) \\ & \geq \pi_{m+p}(u + \epsilon, 0, w) - \pi_{m+p}(u, 0, w) \end{aligned}$$

by (2.19), (2.16), (2.18) and (2.11), (2.15), and (2.12). If  $p = 1$ , then by (2.11) and (2.3b) we may reduce  $w_{m+1}$  to  $\max\{u_{m+1}, w_m\}$  to yield a new policy  $(\tilde{U}, \tilde{Z}, \tilde{W}) \in S(3, r)$  with  $\tilde{m} \geq m+1$  or with  $\tilde{p} > p$ . Now assume

$$(4.4) \quad p \geq 2.$$

Let  $\lambda = \min\{w_{m+1} - w_m, \Delta/p\} > 0$  and define  $(\tilde{U}, \tilde{Z}, \tilde{W})$  by  $\tilde{u}_t = u_t - \lambda$  for  $t = m+1, \dots, m+p-1$ ,  $\tilde{u}_{m+p} = u_{m+p} + (p-1)\lambda$ ,  $\tilde{u}_t = u_t$  for  $t \leq m$  and  $t > m+p$ ;  $\tilde{w}_t = w_t - \lambda$  for  $m+1 \leq t \leq m+p$ ,  $\tilde{w}_t = w_t$  for  $t \leq m$  and  $t > m+p$ ;  $\tilde{z}_t = z_t$  for all  $t$ . To show feasibility of  $(\tilde{U}, \tilde{Z}, \tilde{W})$ , we need only demonstrate  $i_{m+t} \geq t\lambda$  for  $t = 1, 2, \dots, p-1$ , since  $\tilde{z}_t \leq \alpha \tilde{w}_t$  by (4.2). As  $u_{m+p-1} + z_{m+p-1} - d_{m+p-1} = w_r - d_{m+p-1} \geq w_r - d_{m+p} = u_{m+p} + z_{m+p} - d_{m+p} + \Delta$ , we have  $i_{m+p-1} \geq \Delta \geq (p-1)\lambda$ . Let  $j$  be such that  $d_{m+j-1} < w_r \leq d_{m+j}$ , then  $i_{m+j} \geq i_{m+j+1} \geq \dots \geq i_{m+p-1} \geq (p-1)\lambda$ . As  $w_{m+t} - d_{m+t} \geq w_{m+t+1} - d_{m+t+1}$  for  $t = 1, 2, \dots, j$ ,  $i_m \geq 0$ , and  $i_{m+j} \geq (p-1)\lambda$ , we have  $i_{m+t} \geq \frac{t}{j}(p-1)\lambda \geq t\lambda$  for  $t = 1, 2, \dots, j$ .

$$\begin{aligned} & \sum_{t=m+1}^{m+p-1} [\pi_t(w_r, 0, w_r) - \pi_t(w_r - \lambda, 0, w_r - \lambda)] + \pi_{m+p}(u_{m+p}, 0, w_r) \\ & - \pi_{m+p}(u_{m+p} + (p-1)\lambda, 0, w_r - \lambda) \geq \sum_{t=1}^{p-1} \{[\pi_{m+p}(w_r, 0, w_r) - \pi_{m+p}(w_r - \lambda, 0, w_r - \lambda)] \\ & - [\pi_{m+p}(u_{m+p} + t\lambda, 0, w_r - \lambda) - \pi_{m+p}(u_{m+p} + (t-1)\lambda, 0, w_r - \lambda)]\} \geq 0. \end{aligned}$$

The first inequality follows from (2.11), (2.19), and adding and subtracting

$$\sum_{t=1}^{p-2} \pi_{m+p}(u_{m+p} + t\lambda, 0, w_r - \lambda);$$

the second follows from (2.15). Now using the above,  $\lambda \leq w_{m+1} - w_m$  and (2.3b), and  $\tilde{i}_t \leq i_t$  and (2.6) we have  $Q(\tilde{U}, \tilde{Z}, \tilde{W}) \leq Q(U, Z, W)$ . If  $\lambda = w_{m+1} - w_m$ , then  $\tilde{m} > m$ . If  $\lambda = \Delta/p$ , then  $\tilde{p} > p$ . In either case we continue these constructions, mutatis mutandis, until  $\tilde{p} = r$ , in which case  $(\tilde{U}, \tilde{Z}, \tilde{W}) \in S(4, r)$ , or until  $\tilde{m} = r - 1$ . In this case, we have  $\tilde{w}_{r-1} < \tilde{w}_r > \tilde{w}_{r+1}$ ,  $\tilde{z}_r = \tilde{z}_{r+1} = 0$ , and  $\tilde{u}_r < \tilde{w}_r$ . We can now lower  $\tilde{w}_r$  to  $\max\{\tilde{u}_r, \tilde{w}_{r-1}\}$ . If  $\tilde{u}_r \geq \tilde{w}_{r-1}$ , then the new policy, say  $(U', Z', W')$ , is in  $S(4, r)$ . If  $\tilde{w}_{r-1} > \tilde{u}_r$ , then  $m' \leq m - 1$ . Hence,  $m - 1$  more applications of the above yields a policy, say  $(U'', Z'', W'')$ , in  $S(4, r)$  or  $w''_0 = w''_1 = \dots = w''_r > w''_{r+1}$  and  $z''_{r+1} = 0$ . Hence by Lemma 4, there is a policy  $(\bar{U}, \bar{Z}, \bar{W}) \in S(1, r+1)$  such that

$$Q(\bar{U}, \bar{Z}, \bar{W}) \leq Q(U, Z, W).$$

Q.E.D.

LEMMA 6: There is a policy  $(U', Z', W') \in S(5, r) \cup S(1, r+1)$  such that  $Q(U', Z', W') \leq Q(U, Z, W)$  for all  $(U, Z, W) \in S(4, r) \sim [S(5, r) \cup S(1, r+1)]$ ,  $r = 1, 2, \dots, T - 1$ .

PROOF: Take  $(U, Z, W) \in S(4, r) \sim [S(5, r) \cup S(1, r+1)]$ . Then  $w_r > w_{r+1}$ , and  $z_t > 0$  for some  $t$ ,  $m+1 \leq t \leq r$ . Let  $s = \max\{t: z_t > 0, m+1 \leq t \leq r\}$  and  $\Delta = \min\{w_r - w_{r+1}, z_s\} > 0$ . Let

$$u'_t = u_t \quad \text{for all } t \neq r+1, u'_{r+1} = u_{r+1} + \Delta;$$

$$z'_t = z_t \quad \text{for all } t \neq s, z'_s = z_s - \Delta;$$

$$w'_t = w_t \quad \text{for all } t \neq r+1, w'_{r+1} = \max\{w_{r+1}, u_{r+1} + \Delta\}.$$

If  $d_s < w_r$  let  $j = r$ , else let  $j = \min\{t: s \leq t < r, d_t \geq w_r\}$ . Since  $d_s \leq d_{s+1} \leq \dots \leq d_{r+1}$ , we have that  $u_s = u_{s+1} = \dots = u_r = w_r$ ,  $z_{s+1} = \dots = z_{r+1} = 0$ ,  $i_j \geq i_{j+1} \geq \dots \geq i_r \geq w_r - w_{r+1} \geq \Delta$  and  $i_j \geq i_{j-1} \geq \dots \geq i_s \geq z_s \geq \Delta$ . Thus  $(U', Z', W')$  is feasible. From (2.6) and (2.3b) (using  $w_r = w'_r \geq w'_{r+1} \geq w_{r+1}$ ), to show  $Q(U', Z', W') \leq Q(U, Z, W)$  it suffices to demonstrate that

$$(4.5) \quad \pi_S(u_s, z_s, w_r) - \pi_S(u_s, z_s - \Delta, w_r) \geq \pi_{r+1}(u_{r+1} + \Delta, z_{r+1}, w'_{r+1}) - \pi_{r+1}(u_{r+1}, z_{r+1}, w_{r+1}).$$

By (2.16), given  $u_s = w_r$ , (2.18), and (2.19), we have

$$(4.6) \quad \pi_S(w_r, z_s, w_r) - \pi_S(w_r, z_s - \Delta, w_r) \geq \pi_{r+1}(w_r + \Delta, 0, w_r + \Delta) - \pi_{r+1}(w_r, 0, w_r).$$

If  $w'_{r+1} = w_{r+1} \geq u_{r+1} + \Delta$ , then by successive application of (2.13), (2.11), and (2.15) to (4.6), we have (4.5). If  $w'_{r+1} = u_{r+1} + \Delta \geq w_{r+1}$ , then by successive application of (2.13) and (2.11) to

(4.6), we have (4.5). If  $\Delta = w_r - w_{r+1}$ , then  $(U', Z', W') \in S(1, r+1)$ . If  $\Delta \neq w_r - w_{r+1}$ , then  $s' \leq s-1$ . Hence, after at most  $s-m$  such constructions, we have  $(U', Z', W') \in S(5, r) \cup S(1, r+1)$ . Q.E.D.

COROLLARY 7: If  $(U, Z, W) \in S$ , then there is a policy  $(U', Z', W')$  such that

$$(U', Z', W') \in S(5, r) \cup S(1, T)$$

and

$$Q(U', Z', W') \leq Q(U, Z, W).$$

## 5. PROOF OF THEOREM 1

Part (i) ( $w_t \leq w_{t+1}$  for  $t=0, 1, \dots, T-1$ )

Take  $(\tilde{U}, \tilde{Z}, \tilde{W}) \in S$ , then by Corollary 7 there is a policy  $(U, Z, W) \in S(5, \tilde{r}) \cup S(1, T)$  such that  $Q(U, Z, W) \leq Q(\tilde{U}, \tilde{Z}, \tilde{W})$ . Suppose  $(U, Z, W) \notin S(1, T)$ , then there is an  $r \geq \tilde{r}$  and an  $0 \leq m \leq r-1$  such that  $(U, Z, W) \in S(5, r, m)$  and  $w_r > w_{r+1}$ .

Let

$$\lambda_1 = \min \left\{ \frac{w_r - w_{r+1}}{r - m + 1}, w_{m+1} - w_m \right\} > 0$$

$$u_t^1 = \begin{cases} u_t - \lambda_1 & \text{for } t = m+1, \dots, r \\ u_{r+1} + (r+m)\lambda_1 & \text{for } t = r+1 \\ u_t & \text{for all other } t \end{cases}$$

$$w_t^1 = \begin{cases} w_t - \lambda_1 & \text{for } t = m+1, \dots, r \\ \max(w_{r+1}, u_{r+1}^1) & \text{for } t = r+1 \\ w_t & \text{for all other } t \end{cases}$$

$$z_t^1 = z_t \quad \text{for all } t.$$

The feasibility of  $(U^1, Z^1, W^1)$  follows from the argument used to show feasibility of  $(\tilde{U}, \tilde{Z}, \tilde{W})$  in Lemma 5. Since  $w_t^1 = w_t$  for all  $t = m+1, \dots, r+1$ ,  $w_m \leq w_{m+1}^1 = \dots = w_r^1 \leq w_r$ ,  $w_r^1 \geq w_{r+1}^1$ , and  $i_t^1 \leq i_t$  for all  $t$ , we have from Lemma 1 (with  $j = m$ ,  $k = r$ ) (2.16), and (2.19) that  $Q(U^1, Z^1, W^1) \leq Q(U, Z, W)$  if

$$P \equiv \sum_{t=m+1}^r [\pi_{r+1}(w_r, 0, w_r) - \pi_{r+1}(w_r - \lambda_1, 0, w_r - \lambda_1)] \\ - \pi_{r+1}(u_{r+1}^1, 0, w_{r+1}^1) + \pi_{r+1}(u_{r+1}, 0, w_{r+1}) \geq 0.$$

If  $w_{r+1}^1 = w_{r+1} \geq u_{r+1} + (r-m)\lambda_1$ , then adding and subtracting

$$\sum_{t=1}^{r-m-1} \pi_{r+1}(u_{r+1} + t\lambda_1, 0, w_{r+1})$$

we have

$$\begin{aligned} P &= (r-m) \sum_{t=1}^{r-m} \{ [\pi_{r+1}(w_r, 0, w_r) - \pi_{r+1}(w_r - \lambda_1, 0, w_r - \lambda_1)] \\ &\quad - [\pi_{r+1}(u_{r+1} + t\lambda_1, 0, w_{r+1}) - \pi_{r+1}(u_{r+1} + (t-1)\lambda_1, 0, w_{r+1})] \} \\ &\geq (r-m) \{ [\pi_{r+1}(w_r, 0, w_r) - \pi_{r+1}(w_r - \lambda_1, 0, w_r - \lambda_1)] \\ &\quad - [\pi_{r+1}(w_{r+1}, 0, w_{r+1}) - \pi_{r+1}(w_{r+1} - \lambda_1, 0, w_{r+1})] \} \geq 0. \end{aligned}$$

The first inequality follows from (2.15), while the second follows from (2.11) and then (2.13), with  $w = w_{r+1} - \lambda_1$ ,  $\epsilon = \lambda_1$  and  $\beta = w_r - w_{r+1}$ . If  $w_{r+1}^1 = u_{r+1} + (r-m)\lambda_1 > w_{r+1}$ , then  $\pi_{r+1}(u_{r+1}, 0, w_{r+1}) \geq \pi_{r+1}(u_{r+1}, 0, u_{r+1})$  by (2.11), and hence, by adding and subtracting

$$\sum_{t=1}^{r-m-1} \pi_{r+1}(u_{r+1} + t\lambda_1, 0, u_{r+1} + t\lambda_1),$$

we have that  $P \geq 0$  from (2.13) and  $w_r > u_{r+1} + (r-m)\lambda_1$ .

We now recursively define an infinite sequence of policies  $(U^n, Z^n, W^n)$  for  $n=2, 3, \dots$  by

$$\begin{aligned} \lambda_n &= \min \left\{ \frac{w_r^{n-1} - w_{r+1}^{n-1}}{r-m+1}, w_{m+1}^{n-1} - w_m^{n-1} \right\}, \\ u_t^n &= \begin{cases} u_t^{n-1} - \lambda_n & \text{for } t = m+1, \dots, r \\ u_{r+1}^{n-1} + (r-m)\lambda_n & \text{for } t = r+1 \\ u_t^{n-1} & \text{for all other } t \end{cases} \\ w_t^n &= \begin{cases} w_t^{n-1} & \text{for } t = m+1, \dots, r \\ \max \{ w_{r+1}^{n-1}, u_{r+1}^n \} & \text{for } t = r+1 \\ w_t^{n-1} & \text{for all other } t \end{cases} \end{aligned}$$

$$z_t^n = z_t^{n-1} \quad \text{for all } t.$$

Employing the argument above, we have, for all  $n$ , that  $Q(U^{n+1}, Z^{n+1}, W^{n+1}) \leq Q(U^n, Z^n, W^n)$  and that

$$(U^n, Z^n, W^n) \in \bigcup_{j=0}^m S(5, r, j) \cup S(1, r+1) = A$$

implies  $(U^{n+1}, Z^{n+1}, W^{n+1}) \in A$ .

CASE 1: There is an  $n$ , say  $N$ , such that  $w_m^N = w_{m+1}^N$ . Let  $m^N$  denote the  $m$  corresponding to policy  $(U^N, Z^N, W^N)$ .

In this case  $m^N < m$ . Now letting  $(U^N, Z^N, W^N)$  replace  $(U, Z, W)$  in all of the above and continuing the constructions, we arrive at either Case 2 below, or, after at most  $m$  such replacement at a policy  $(\tilde{U}, \tilde{Z}, \tilde{W})$ , such that  $\tilde{w}_0 = \tilde{w}_1 = \dots = \tilde{w}_r > \tilde{w}_{r+1}$ , in which case, we have by Lemma 4 that there is a lower cost policy in  $S(1, r+1)$ .

CASE 2: There is no  $n$  such that  $w_m^n = w_{m+1}^n$ .  
Then

$$w_t^\infty \equiv \lim_{n \rightarrow \infty} w_r^n = \lim_{n \rightarrow \infty} w_{r+1}^n \equiv w_{r+1}^\infty \quad \text{for } t = m+1, \dots, r$$

since  $w_t^n = w_r^n$  for  $t = m+1, \dots, r$  and all  $n$  and

$$0 \leq w_r^n - w_{r+1}^n \leq \left( \frac{1}{r-m+1} \right)^{n-1} (w_r^1 - w_{r+1}^1), \quad w_r^{n+1} < w_r^n,$$

and

$$w_{r+1}^{n+1} \geq w_{r+1}^n \quad \text{for all } n.$$

Consequently,  $u_t^\infty \equiv \lim_{n \rightarrow \infty} u_r^n = \lim_{n \rightarrow \infty} w_r^n = w_r^\infty$  for  $t = m+1, \dots, r$ , since  $u_t^n = w_t^n$  for  $t = m+1, \dots, r$ ,  $n = 1, 2, \dots$ .  $u_{r+1}^\infty \equiv \lim_{n \rightarrow \infty} u_{r+1}^n$  as  $u_{r+1}^{n+1} \geq u_{r+1}^n$  and  $u_{r+1}^n \leq w_{r+1}^n \leq w_{r+1}^\infty$  for all  $n$ . Thus setting  $z_t^\infty = z_t$  for all  $t$  and  $u_t^\infty = u_t$ ,  $w_t^\infty = w_t$  for  $t \neq m+1, \dots, r+1$ , we have that  $(U^\infty, Z^\infty, W^\infty) \in S(1, r+1)$ .

In either case, we arrive at a policy in  $S(1, r+1)$ . Thus repeating the entire process at most  $T-r$  times, we find a policy, say  $(U', Z', W')$ , in  $S(1, T)$ , and so satisfying (i), such that  $Q(U', Z', W') \leq Q(U, Z, W)$ . Q.E.D.

Part (ii) ( $u_t = w_t$  for  $t = 1, 2, \dots, T$ )

Let  $(\tilde{U}, \tilde{Z}, \tilde{W})$  be any policy. If  $\tilde{u}_t \geq \tilde{w}_t$  for  $1 \leq t \leq T$ , then let  $\tilde{r} = T+1$ . Otherwise,  $\tilde{r} \equiv \min\{t: \tilde{u}_t < \tilde{w}_t, 1 \leq t \leq T\}$ .

\*If  $w_r - w_{r+1} = \eta$  and  $w_{r+1} - u_{r+1} > \eta$ , then  $w_{r+1}^\infty = w_{r+1}$ ,  $u_{r+1}^\infty = w_{r+1}^\infty$ , and  $u_t^\infty = w_t = w_{r+1}^\infty$  for  $t = m+1, \dots, r$ .



Let  $P(r) \equiv \{(\tilde{U}, \tilde{Z}, \tilde{W}) \in S : (\tilde{U}, \tilde{Z}, \tilde{W}) \text{ satisfies (i) and } \tilde{r} \geq r\}$ , for  $r = 1, 2, \dots, T$ . Take  $(U, Z, W) \in P(r) \sim P(r+1)$ .

CASE 1:  $w_0 = w_1 = \dots = w_r$ . Using  $(U, Z, W) \in S$ , (3.1), (3.2), and  $u_r < w_r$ , we have

$$\sum_{t=1}^r z_t \geq w_r - u_r.$$

Let  $w'_t = w_t$  for all  $t$ ;  $u'_t = u_t$  for all  $t \neq r$ ,  $u'_r = w_r$ ;

$$z'_t = z_t, t > r, z'_t = \max \left\{ z_t - \left[ w_r - u_r - \sum_{i=t+1}^r (z_i - z'_i) \right], 0 \right\}, \quad t \leq r.$$

By applying (2.19), (2.16), (2.18), (2.11), (2.15), and (2.12), we have

$$\begin{aligned} \pi_t(w, z, w) - \pi_t(w, z - \epsilon_t, w) &\geq \pi_r(w, z, w) - \pi_r(w, z - \epsilon_t, w) \\ &\geq \pi_r(w, \epsilon_t, w) - \pi_r(w, 0, w) \geq \pi_r(w + \epsilon_t, 0, w + \epsilon_t) - \pi_r(w, 0, w) \\ &\geq \pi_r(w + \epsilon_t, 0, w + \epsilon_t) - \pi_r(w, 0, w + \epsilon_t) \geq \pi_r(u_t + \epsilon_t, 0, w + \epsilon_t) \\ &\quad - \pi_r(u_t, 0, w + \epsilon_t) \geq \pi_r(u_t + \epsilon_t, 0, w) - \pi_r(u_t, 0, w), \end{aligned}$$

where  $\epsilon_r = z_r - z'_r$ ,  $\epsilon_t = z_t - z'_t$ ,  $w = w_r$ , and

$$u_t = u_r + \sum_{i=t+1}^r \epsilon_i,$$

for  $t = 1, 2, \dots, r-1$ . Using (2.17) for  $t = r$ , the above, and (2.6) we have  $Q(U', Z', W') \leq Q(U, Z, W)$  and  $Q(U', Z', W') \in P(r+1)$ .

CASE 2:  $w_m < w_{m+1} = \dots = w_r$ ,  $0 \leq m < r$ . As shown in Case 1 above, we may assume  $z_t = 0$  for  $t = m+1, \dots, r$ . Let  $\lambda = \min\{w_{m+1} - w_m, (w_r - u_r)/(r - m)\}$

$$z'_t = z_t \quad \text{for all } t;$$

$$u'_t = u_t, w'_t = w_t \quad \text{for } t \neq m+1, \dots, r;$$

$$u'_t = w'_t = w_r - \lambda, \quad m+1 \leq t \leq r-1,$$

$$u'_r = u_r + (r - m - 1)\lambda, \quad w'_r = w_r - \lambda.$$

The argument of Lemma 5 shows that  $(U', Z', W') \in S$  and  $Q(U', Z', W') \leq Q(U, Z, W)$ . If  $\lambda \neq (w_r - u_r)/(r - m)$ , then  $m' < m$ . In which case at most  $m - 1$  repetitions are needed until either  $\lambda = (w_r - u_r)/(r - m)$  or  $w_0 = w_1 = \dots = w_r$ . In the latter case, we have (from Case 1 above) the existence of  $(U'', Z'', W'') \in P(r+1)$  with  $Q(U'', Z'', W'') \leq Q(U, Z, W)$ . If  $\lambda = (w_r - u_r)/(r - m)$ , then  $(U', Z', W') \in P(r+1)$ .

Thus, if  $(U, Z, W) \in P(r) \sim P(r+1)$ , we can find a policy  $(U', Z', W') \in P(r+1)$  such that  $Q(U', Z', W') \leq Q(U, Z, W)$ . Q.E.D.

Part (iii) and (iv) ( $u_t \leq u_{t+1}$  and  $u_t + z_t \leq u_{t+1} + z_{t+1} \leq d_t$  for  $t = 1, 2, \dots, T-1$ )

(iii) follows immediately from parts (i) and (ii). Let  $(\tilde{U}, \tilde{Z}, \tilde{W})$  be any policy. If  $\tilde{u}_t + \tilde{z}_t \leq \tilde{u}_{t+1} + \tilde{z}_{t+1}$  for  $1 \leq t \leq T-1$  let  $\tilde{r} = T$ . Otherwise  $\tilde{r} = \min\{t : \tilde{u}_t + \tilde{z}_t > u_{t+1} + z_{t+1}, 1 \leq t \leq T-1\}$ . Let  $P(r) \equiv \{(\tilde{U}, \tilde{Z}, \tilde{W}) \in S : (\tilde{U}, \tilde{Z}, \tilde{W}) \text{ satisfies (i)-(iii) and } \tilde{r} \geq r\}$ , for  $r = 1, 2, \dots, T-1$ . Take  $(U, Z, W) \in P(r) \sim P(r+1)$ , so that

$$(5.1) \quad u_r + z_r > u_{r+1} + z_{r+1}.$$

From (iii), (2.2), and (5.1), we have  $z_r \geq z_r - z_{r+1} > u_{r+1} - u_r \geq 0$ . Also from (3.2)

$$i_r \geq \Delta \equiv \frac{1}{2} (u_r + z_r - u_{r+1} - z_{r+1}) > 0. \text{ Let}$$

$$u'_t = u_t, \quad w'_t = w_t \quad \text{for all } t,$$

$$z'_t = z_t \quad \text{for all } t \neq r, r+1;$$

$$z'_r = z_r - \Delta, \quad z'_{r+1} = z_{r+1} + \Delta.$$

As  $\Delta \leq \frac{1}{2} (z_r - z_{r+1})$ , we have  $z' \geq 0$  and  $z'_{r+1} \leq z_{r+1} + \frac{1}{2} (z_r - z_{r+1}) < z_r \leq \alpha w_r \leq \alpha w_{r+1} = \alpha w'_{r+1}$ . Hence  $(U', Z', W') \in P(r+1)$ . Thus, it remains only to show that  $Q(U', Z', W') \leq Q(U, Z, W)$ . Let  $\epsilon = w_{r+1} - w_r = u_{r+1} - u_r$ , then  $z_r = 2\Delta + \epsilon + z_{r+1}$  and so by (2.19) and (2.16), and (A1)

$$\begin{aligned} & \pi_r(w_r, z_r, w_r) - \pi_r(w_r, z_r - \Delta, w_r) \\ & \geq \pi_{r+1}(w_r, z_{r+1} + \epsilon + \Delta, w_r) - \pi_{r+1}(w_r, z_{r+1} + \epsilon, w_r) \\ & \geq \pi_{r+1}(w_r + \epsilon, z_{r+1} + \Delta, w_r + \epsilon) - \pi_{r+1}(w_r + \epsilon, z_{r+1}, w_r + \epsilon). \end{aligned}$$

Q.E.D.

## 6. ASYMPTOTIC PROPERTIES

We now turn attention to characterizing the form of optimal policies as the horizon lengthens. Our orientation is to consider questions of long range planning, and in particular, we seek information about the values of the instrumental variables at a period far in the future. We extend the notation in the obvious way to  $u_t(T)$ ,  $z_t(T)$ , and  $w_t(T)$ .

Let  $\lim_{t \rightarrow \infty} d_t \equiv B \leq +\infty$ . If  $B = +\infty$ , then feasibility considerations alone imply that,

given  $M \geq 0$ , there exists a period  $n$  such that

$$w_n(T) \geq u_n(T) \geq M \quad \text{for } T \geq n.$$

Suppose instead that  $B < \infty$ . We might expect that far in the future  $u_t(T)$  is approximately  $B$  and  $z_t(T)$  is arbitrarily small. To ensure this result, we must add an assumption on  $\pi_t(u, z, w)$ .

We assume that for each specified pair  $(\epsilon, Y)$ ,  $\epsilon > 0$  and  $Y$  a nonnegative integer, there exists a nonnegative integer  $Z(\epsilon, Y)$  such that

$$(6.1) \quad \sum_{t=Y}^{Z(\epsilon, Y)} P_t(\epsilon) > [g_Y + f_{Z(\epsilon, Y)}] \epsilon,$$

where

$$(6.2) \quad P_t(\epsilon) = \inf \{ \pi_t(w, z, w) - \pi_t(w + \epsilon, z - \epsilon, w + \epsilon) : \epsilon/\alpha \leq w \leq B - \epsilon, \epsilon \leq z \leq \alpha w \}.$$

If

$$(6.3) \quad \pi_t(\cdot, \cdot, \cdot) = \beta^{t-1} \pi_0(\cdot, \cdot, \cdot), f_t = \beta^{t-1} f, \text{ and } g_t = \beta^{t-1} g \quad \text{for } t = 1, 2, \dots,$$

then (6.1) is equivalent to

$$(6.4) \quad P_0(\epsilon)/(1 - \beta) > g\epsilon, \quad \epsilon > 0;$$

thus (6.1) rules out some examples with discounting. Suppose, in addition to (6.3) that  $P(\epsilon)$  is bounded away from zero, as is the case when production costs are linear, then there exists a number  $\beta_0$  ( $0 \leq \beta_0 < 1$ ) such that (6.1) holds if, and only if,  $\beta \in (\beta_0, 1)$ .

**THEOREM 2:** If  $B < \infty$  and (6.1) holds, then given  $0 < \delta < B$ , there exists a period  $n(\delta)$  such that, when  $T \geq n(\delta)$ , we have for some optimal policy that

$$(6.5) \quad B - \delta \leq u_{n(\delta)}(T) \leq \dots \leq u_T(T) \leq B$$

$$(6.6) \quad 0 \leq z_t(T) \leq \delta \quad \text{for } n(\delta) \leq t \leq T$$

$$(6.7) \quad B - \delta \leq w_{n(\delta)}(T) \leq \dots \leq w_T(T).$$

**PROOF:** Let  $\langle (U(T), Z(T), W(T)) \rangle_{T=1}^{\infty}$  be a sequence of optimal policies satisfying

Theorem 1, and let  $0 < \delta < B$  be given. Choose  $\bar{t}$  such that  $d_{\bar{t}} > B - \frac{\delta}{4}$ . Then, since

$i_{\bar{t}} \leq B\bar{t} + i_0$  there exists a period  $N$  such that  $\frac{\delta}{4} N \geq i_{\bar{t}}$ . Let  $t^* = N + \bar{t}$  and take  $T \geq t^*$ . If

$u_{t^*}(T) + z_{t^*}(T) < B - \frac{\delta}{2}$ , then by (iv) of Theorem 1 we have:

$$i_{t^*}(T) = i_{\bar{T}} + \sum_{t=\bar{T}+1}^{t^*} [u_t(T) + z_t(T) - d_t] < i_{\bar{T}} + N\left(B - \frac{\delta}{2}\right) - N\left(B - \frac{\delta}{4}\right) < 0.$$

Hence by (2.5)

$$(6.8) \quad u_{t^*}(T) + z_{t^*}(T) \geq B - \frac{\delta}{2} \quad \text{for all } T \geq t^*.$$

Let  $n(\delta) = Z\left(\frac{\delta}{2}, t^*\right)$  and fix  $T \geq n(\delta)$  (we now suppress the  $T$ ).

Suppose  $u_{n(\delta)} < B - \delta$ , then  $z_t > B - \frac{\delta}{2} - (B - \delta) = \frac{\delta}{2}$  for  $t^* \leq t \leq n(\delta)$  by (6.8) and (iii) of Theorem 1. Now let

$$u'_t = u_t, z'_t = z_t, w'_t = w_t \quad \text{for } t < t^*, n(\delta) < t \leq T,$$

$$u'_t = u_t + \frac{\delta}{2}, z'_t = z_t - \frac{\delta}{2}, w'_t = w_t + \frac{\delta}{2} \quad \text{for } t^* \leq t \leq n(\delta).$$

Thus by (6.1)

$$\begin{aligned} Q(U, Z, W) - Q(U', Z', W') &= \sum_{t=t^*}^{n(\delta)} \left\{ \pi_t(w_t, z_t, w_t) - \pi_t\left(w_t + \frac{\delta}{2}, z_t - \frac{\delta}{2}, w_t + \frac{\delta}{2}\right) \right\} \\ &\quad - [g_{t^*} + f_{n(\delta)}] \frac{\delta}{2} \geq \sum_{t=t^*}^{n(\delta)} P_t\left(\frac{\delta}{2}\right) - [g_{t^*} + f_{n(\delta)}] \frac{\delta}{2} > 0, \end{aligned}$$

contradicting the optimality of  $(U, Z, W)$ .  $u_T \leq B$  and (6.6) follow from (2.6) and (2.11). (6.7) follows from (2.1). Q.E.D.

## REFERENCES

- [1] Lippman, S. A., A. J. Rolfe, H. M. Wagner, and J. S. C. Yuan, "Optimal Production Scheduling and Employment Smoothing with Deterministic Demands," *Management Science*, 14, 127-158 (Nov. 1967).
- [2] Lippman, S. A., A. J. Rolfe, H. M. Wagner, and J. S. C. Yuan, "Algorithms for Optimal Production Scheduling and Employment Smoothing," *Operations Research*, 15, 1011-1029 (Nov.-Dec. 1967).

\* \* \*

A. S. Goldman

*General Electric Company*  
*TEMPO*  
*Santa Barbara, California*

## SUMMARY

Although industry is expected to design hardware to fit into a general support system and to be capable of arguing life-cycle system costs, adequate information has not been available on the support system in terms of policies and operating decision rules. Policies and operating decisions by users dominate engineering design decisions in determining life-cycle support costs. The relative effect of each of these decision areas on support costs has yet to be resolved empirically. Without an understanding of the sensitivity of support costs to alternative designs, capability is limited in design improvement and support of end items. Life-cycle costing of analysis under cost-effectiveness and the maintainability of integrated logistics support is open to question.

## BACKGROUND

Depending on the fiscal year and what is interpreted as "cost of support," it is generally agreed that at least 20-30 percent of the total defense dollar goes to the support of existing equipments; and that, depending on the system or equipment and time period considered, life-time support costs often dominate system life cycle costs by an order of at least ten times the original cost of the end unit.\* These are the two foremost reasons given for the current interest in support problems, particularly as support costs are involved in such concepts as maintainability, integrated logistics support, life cycle costs, and cost-effectiveness analysis.

The fact that support costs are high requires an understanding of what decisions go into support and the requirements that result so that approaches to support cost estimation can be agreed upon. While the concept of costing out of support seems simple enough on the surface, there are a number of pitfalls which can lead to misleading decisions in system design and procurement.

The objective of this paper is to attempt to outline the difficulties in costing out support. These problems dominate the discussion:

---

\*This relationship of support cost to acquisition cost naturally varies widely depending on the system. Other than parametric studies, limited empirical experience exists for evaluating this relationship by weapon system. This in a sense is what this paper is about: To specify the reasons for our inability to identify empirically support costs of a weapon system. There are planning factors that permit estimates to be made of the support cost relationship, but these are of a gross nature for aid in high level decisions in contrast to decisions involved with engineering design within a weapon system. Moreover, certain of these factors require further justification.



- The influence of management decisions of users (divided into policy and operating decisions) on support costs.
- The partitioning of the support system in which depot level support (supply, overhaul, and repair) is guided by criteria which are only indirectly related to operational support and end unit effectiveness.
- Support system stocks of resources are influenced by initial investment in support, which tends to be insensitive to many types of individual engineering design decisions.

The intent is not to present ways of approaching these costing difficulties. Once the problem areas are outlined, however, it becomes evident that problems which are currently set aside explicitly or implicitly by assumption, can be handled from a practical point of view by a meeting of the minds between analyst and reviewer in the course of evaluating cost differences due to design differences.

The primary issue that evolves in the discussion that follows is that support system management decisions over the life cycle overshadow most individual design decisions in the effect on relevant\* life cycle support costs. Unless we are able to distinguish between costs due to management decisions and those due to design differences, and unless we are able to delineate what is meant by a weapon system and its associated support resources, different designs could well have similar lifetime support costs; conversely, similar designs could have different lifetime support costs.

One final introductory point: If concern is with the relative support costs for two or more designs having the same performance capability, this can reasonably be handled since it is then assumed that the key determining factor influencing support costs will be failure rates—everything else that can affect support costs is presumed constant.† But this problem is different from one in which absolute life cycle costs are at issue, i.e., where cost and effectiveness are involved in evaluating a weapon system as it fits into the total force structure (where performance and mission objectives differ), or where life cost estimation is required for budgeting and financial purposes (as in the LMI effort in Integrated Logistics Planning‡).

## USER INFLUENCE ON SUPPORT COSTS

In this paper, we will distinguish between two kinds of decisions that determine support costs: management (of the user) and engineering (of the contractor).§ An additional factor common to both decision areas is the discrete nature of support resources which tends to make differences in support costs insensitive to many kinds of individual design alternatives.

\*In contrast to irrelevant costs--where relevant is synonymous with "out-of-pocket" and irrelevant with "sunk" costs.

† This statement may be acceptable for practical reasons, but strictly speaking, assumes that systems using existing support assets should have more favorable support costs because (a) they make more efficient use of common resources, and (b) a proportion of common and existing assets are sunk costs, not charged against the system.

‡ Logistics Management Institute, DOD Systems and Equipment Integrated Logistics Support Planning Guide, LMI Task 66-15, December 1967.

§ Although we are concerned here with support costing of a given system, note that one area that has had little explicit study involves the impact of performance requirements on support costs. Tradeoff studies between performance capability (related to the mission and determining end unit complexity and reliability) and availability (performance over time) could prove of interest so far as identifying a key influence on support costs and, thereby, system choice.

In order to develop this thinking, we take the point of view that, subject to maintenance policies on repair and overhaul, the end objective of the support system is to provide for required spares and repair parts. Thus, for any predetermined maintenance policy, the argument centers on inventory management and the allocation of spares and repair parts as the determinant of the effectiveness of the support system and as the key cost variable in support which reflects operating decisions.

#### MANAGEMENT DECISIONS BY USERS

"Management decisions" refers to those choices of actions leading to future events desired by the decision maker at the time the action is taken. Choice may result from quantitative (either deterministic or probabilistic) analysis or from intuitive and subjective reasoning (or some combination which we can omit for purposes of this paper). The operation of support systems is characteristically the former in which recurring decisions are made day-by-day within the existing support system; the latter relates to "policy," one-shot or infrequently made decisions which either constrain the sequence of operating decisions or structure the system in which operating decisions are made.

#### OPERATING DECISIONS

Recurring decisions in an operating support system have as their objective the maintenance of end units and availability of spares and repair parts by means of continual replenishment of stocks. The replenishment rules for stocks, involving regenerations from repair or from procurement, determine asset levels and requirements in the support of existing primary systems; these rules we group under quantitative areas of management.

From an operational profile of requirements, the support system reacts to hardware component characteristics (design decisions determining complexity of components which in turn relate to item lead times for procurement, failure rates, redundancy, and so on) according to procedures and rules to maintain a pipeline of material. Subject to preventive maintenance policies, replenishment rules and procedures in the repair or procurement process include:

- Screening decisions involving repair or discard.
- Economic lot size for repair, procurement, or transportation.
- When to repair, reorder, or transportation.
- The size of pools of stocks for safety either in the repair process or at stocking points.

#### USER POLICY DECISIONS

Over and above these operating decision rules are the subjective rules, policies which dictate the framework of environment within which operating decisions are made. Policy rules are largely courses of action selected from among alternatives which are not subject to quantitative analysis, but which serve to guide and bound present and future operating decisions. The evaluation of a policy decision, constrained as it usually is by complexity of the support system and lack of an empirical base, often results in a final choice among alternative policies which is largely judgmental. Examples of policy areas that affect support costs:

- The extent that any new system can or should draw upon common facilities and stocks of the general support system.

- Degree that operational (organizational) levels will be self-sustaining. Level of maintenance and repair provided at operational, intermediate, and depot levels; and at each level, the repair/discard policy to be followed.\*

- Man loading allocation at each echelon and their support (personnel and facilities).
- Protection against shortages: protection against shortages at operational level or against shortages of common support assets at depot level which can degrade an operating system in the field. System protection as contrasted with protection at an echelon.
- Range of stocks to be carried at echelons.
- Extent that premium transportation ties in with stocking policies.
- Time period between overhaul of primary units.

### THE PARTITIONING PROBLEM

By "partitioning of the system" we mean that local decision rules at depot level (where several weapon programs are supported) are based upon usage and management of items, many of which are used by several weapons, without regard to individual weapon system requirements; that criteria at the depot such as "fill rate" (percent of items demanded that are satisfied) are only indirectly related to end user effectiveness; that depot level concern is with budgets and costs, stock turnover in terms of asset levels, and sales. At the operational and intermediate levels the objective is, e.g., weapon or program readiness and sustaining of allowances. The dollar constraint seldom is a dominant issue at operating sites.

It is this split between environment and objectives of general support management and end unit weapon system management that can result in stocks at the depot level varying dollar-wise independently of measures of system effectiveness. This is, in part, the stock-flow issue described below, but it also relates to the common asset problem and what is meant by a weapon system program including end unit support.

A cost-effectiveness analysis, where support costs are the issue, begs for identity of costs peculiar to a weapon and common to that weapon and others. From a life-cycle cost point of view, no single system may be viewed in its requirement for annual support and replenishment without recognition that each weapon system competes for common assets with other weapons being supported from the same "common bin"; and that each weapon competes for funds in annual support for items common as well as peculiar to that weapon.† (To repeat what was said initially: costs of particular concern, common or unique, include only out-of-pocket and not sunk costs, which naturally favors systems using existing support assets.)

Ideally, to cost-out a weapon system program, we should be able to bound that portion of a general support system allocated to it (much as in overhead in private industry), and the

\*Subject to deployment requirements, i.e., repair capability of intermediate maintenance is determined by the distance from depot level of supply and repair.

†Recent action at DoD is directed toward identifying usage of items by weapon. NAVAIR (ASO) is now orienting toward management of replenishment spares and parts by weapon. Note that support resources may be divided into three groups: (1) those peculiar to the weapon; (2) common to the weapon system and others; (3) unique to others. Demands for these resources, however, are divided into four areas: common and peculiar to the weapon and common and peculiar to all other weapons.



specific operating decision rules that will determine annual replenishment of both items unique and common to that weapon system program. The problem is that depot level resources are managed not for weapons, as such, but rather on items and stocks that can build up without direct relationship to effectiveness or specific weapon system performance at the operational level.\*

We may speak of weapon program management, but even the most major weapon systems (such as FBM) must be viewed as they fit into depot support, and here is the rub: readiness is an issue at operational levels, and item management is the issue at depot level. Measures of performance for management purposes in the two levels will not necessarily correlate. This is the central issue in determining levels of pipeline stocks to carry at depot level. Based upon a support structure and associated pipeline stocks, which reflect operating decision rules and policies in procurement and repair, an eat-down of depot level stocks can occur within a finite time period without any change in operational level effectiveness of selected systems. This is true not only for stocks directly related to end unit performance, but also for those that are less immediately "vital" or essential to end unit performance.

Conceptually, one can analytically probe depot stock levels to estimate that point that first degrades end unit performance. Above that point, a host of decision rules can exist, particularly affecting economic order quantities and safety levels, which determine support costs, but not end unit performance. In theory, there should be a relation between, e.g., percent protection against shortages at depot and end unit performance and effectiveness for selected weapon programs. But this transformation has as yet no empirical justification. In light of the importance of depot stocks and other resources to cost-effectiveness analysis, this is a major hurdle to overcome in the support costing problem.

## STOCK VERSUS FLOW

A support system is characterized by its capability to utilize the various support system resources in satisfying demands placed upon it by primary units. The key variable resource, subject to the resources established and fairly well fixed in the support system, is spares and repair parts. "Stocks" of spares and repair parts flow between echelons according to decision rules (that determine echelon stocking levels) in response to maintenance and equipment usage requirements. So that, in addition to fixed support resources that exist regardless of operations of primary units, the variable resource in support is in logistics and pipeline assets in support of end units. For this reason, in addition to the importance of logistics costs, we will continue to center our remaining discussion on logistics management to illustrate the difficulty of costing this area of support.

The idea of pipeline stocks within any fiscal period can be visualized in terms of the kinds of depot level requirements used in the DoD replenishment budget process:

- Contingency reserves
- Backorders -- quantities of items on order
- Safety level -- to cover variations in demand and replenishment lead times
- Repair cycle -- recovery of reparables returned for repair

---

\*Communication or transportation systems in support of several weapon systems are particularly difficult to debate in a cost-effectiveness analysis.

- Production lead time — to cover demands between contract award and time of delivery
- Future issues — operating demands.

While this breakdown of pipeline requirements assumes a capability of repair and transportation, it does serve to illustrate the elements that go into determining life cycle stocking requirements and the kinds of management decisions that affect the level of asset requirements. Within this concept of pipeline stocks lies also the subtlety of individual decision rules for hundreds of thousands of secondary or tertiary items — some common to several weapons, some peculiar; some essential to weapon readiness, some not essential — each of which is in a state of "overage" or "underage" relative to expected primary unit requirements within a fiscal period. Expected pipeline requirements, i.e., on the average, is at best a misleading although necessary simplifying conceptual device for management purposes.

Of particular importance, note that the average stock in the pipeline can range widely within a finite time period without relating to primary unit effectiveness.\* Moreover, while pipeline requirements increase for peculiar and common material with additions of new end items, the incremental amount promises to remain an unknown in the foreseeable future. Certainly this will remain an issue until decision rules and their impact on asset levels are agreed upon.† The question here is this: If failure rates and repair rates vary by ten percent in alternative designs, what is the impact on support costs of that weapon when much of this comes out of common bins which support several competitor weapons or when replenishment decisions are a function of funding availability, different pipeline requirements, and user policies on support?

#### INITIAL SET-UP COSTS

A major element of the stock-flow issue is in decisions concerning deployment of primary units at operating sites, the numbers of primary units per site, and the level of maintenance and support at each site as well as in the back-up system. Initial set-up costs, depending on complexity of the system and the use made of existing assets, can run over fifty percent of lifetime support costs. The character of initial buys is largely for system life, i.e., high cost reparables and fixed facilities and test equipments, in contrast to replenishment which is made up of a larger percentage of lower cost consumables.

The importance of initial requirements to system life support costs is largely a matter of policy — the range of initial stocks, the depth of these stocks based on a desired protection against shortages for a fixed time period, the level of repair and associated maintenance manpower and test equipment; and most important, the number of sites and end units per site under a deployment plan wherein greater economies of support are accomplished with larger numbers of end units per site. There is no problem conceptually of prorating these initial costs over lifetime use, but this does not minimize the significance of:

- A deployment plan that spells out numbers of sites and numbers of end units per site and their initial set-up costs;
- Initial support costs as a percentage of lifetime support costs which is largely determined by policy (on preventive maintenance or repair) and the deployment plan;

\*This statement holds for those items essential to primary unit effectiveness as well as those whose essentiality is less direct.

†Realistically, this kind of knowledge will not be available. Expected consumption, exclusive of pipeline requirements, becomes the most meaningful measure of spares and parts cost.



- Risk of cost obsolescence of initial stocks in the life cycle due to modifications;
- The limitations of failure rate projections for new items or for low-usage items;\*
- That lifetime pipeline asset levels are determined by policy decisions on initial set-up requirements and maintenance both at operational and system backup levels;
- An initial maintenance plan concerned with the relative costs of repair at the intermediate and depot levels which is complicated by relative distances between intermediate and depot levels in the deployment plan. Does one plan unknowingly favor repair and therefore that design leaning to repair?

#### FORMAT FOR EVALUATION AND COSTING PURPOSES

In summary, support cost as a decision-making issue is of interest in two areas: (1) during engineering design in the choice between two or more design alternatives, and here relative support costs are the issue; and (2) during force structure type analyses or in budgeting and financial problems in planning for the relevant, out-of-pocket, life cycle costs. To estimate life cycle costs, we have attempted to show that we must be able to distinguish between the contribution to costs by management decisions (of the user) external to or beyond the control of the designer and systems analyst, and those decisions over which the designer has direct control.

Ideally, what sort of information should we have in order to evaluate support costs associated with a design decision or system alternative? In the following, we discuss this information, continuing to draw a distinction between management and engineering decision areas.

An obvious starting point is in classification of equipments in terms of functional use according to primary, secondary, and tertiary roles.† This involves multi-purpose end units, such as machine guns which can be used by infantrymen, or can be mounted on tanks, carriers, and aircraft, and single purpose end units, such as ground-to-air missiles. The classification problem is an arbitrary process that applies to subassemblies of subsystems, such as communication or navigation gear and computers and fire control systems which, while procured for one weapon system, may be used in others.

The point to classification is that the costing of support of an equipment hinges to a large degree on the intended use (or combinations of use). This is, in fact, also part of the common asset problem. Just as there are many purposes for the same item, so are there many uses by different weapon systems and equipments for the same subsystems and assemblies. The common asset problem in support reflects not only upon the end use of an equipment, but also the degree to which an equipment can make use of existing support resources.‡ While we recognize that a system that makes use of existing support system resources can reduce support costs in contrast to a system that must fall back on resources peculiar to that system, the extent of this efficiency and economy remains an unknown.

\*The low-usage item is a phenomenon of inventory management. In major complex systems, some 60 to 70 percent of all stocked items will have one demand per year or less although these may not be the same items from year to year. Usage of these items is unpredictable and not susceptible to quantitative analysis for forecasting purposes within any one year.

†This is related to, but in practice need not be as complex as, engineering analyses of the essentiality of equipments to weapon system effectiveness.

‡The common asset problem is of importance in view of recent decisions by OSD (Comptroller) on increasing use of the Stock Fund and Defense Supply Fund. In fact, the common asset problem affects economies of production as well as support in stocking and repair.

Table 1 illustrates the sort of equipment, mission, and support relationship expressed above.

TABLE 1  
Equipments and Support Related to Missions

Equipment and Support Services	Missions			
	$M_1$	$M_2$	$M_3 \dots$	$M_n$
$E_i$				
$e_{ij}$				
$S_{ij}$				

The  $E_i$  refers to the weapon system program equipments used to perform the  $M_n$ 's (missions); the  $e_{ij}$  refers to the subequipments used in an  $E_i$  in which the subscript "j" relates to other  $E_i$  equipments employing the  $e_{ij}$  subequipment. An  $e_{ij}$  can also become an  $E_{ij}$  depending on the  $M_n$  involved. The  $S_{ij}$  refers to the support resources used to support the  $e_{ij}$ 's and the  $E_i$ 's.

Conceptually, by weighting the  $M_n$ 's, the requirements or expected use of the set of  $E_i$ 's determine, in turn, the  $e_{ij}$ 's and spell out a requirement for the  $S_{ij}$ 's. The total number of equipments and their support can then be priced out depending on the deployment and maintenance plan and number of years involved.

The objective of system costing is to be able to approximate relevant costs associated with a weapon system program — including acquisition costs (research and development where appropriate), initial first year set-up costs, and follow-on annual costs. Essential to costing, therefore, is the identification of end item functional classification by weapon system, or where subassemblies are involved, the association of subassemblies to item use. The concept appears simple, but its practical application and implementation for a mix of varied kinds of weapon subassemblies "across the board" is a complicating feature that we feel can only be resolved by the buyer establishing the basic rules of the game with the contractor by defining  $M_n$ 's and their weights. Any other uses of the  $E_i$ 's or  $e_{ij}$ 's that the contractor can develop is a "plus" factor.

It seems evident from the above that the policy area emphasizes the relation of equipments to functions and to missions. Without this point of view which basically establishes the identification and classification of the weapon system and ties in closely to the objective of effectiveness analysis, an approach to costing would end with competitors coming up with cost estimates that are not comparable. This is equally true in costing for engineering design alternatives where a weapon system or subassembly, which can be used in several functional ways, might be more costly, yet provide more mission capability than a system with less flexibility of use. This thinking of course applies to Services — e.g., the F-111 or F-4. In keeping with

the common asset or standardization concepts, the support costs of a flexible or common weapon system would be less than a highly specialized one, ceteris paribus; we have yet to come to a common agreement as to how much less.

### SUPPORT COSTS, THE $S_{ij}$ 's

Now we refer to the  $S_{ij}$ 's of Table 1, the variable operating costs as contemplated with fixed assets that do not relate to different levels of operation. These costs relate to asset and funding requirements during the life cycle which are a function of operating decision rules, subject to policy decisions. The policy aspect of the  $S_{ij}$ 's is best seen in terms of the initial investment support costs\* imposed by virtue of:

- The range and depth of spares and parts at operational, intermediate and depot levels subject to a requirement (by policy) on protection against shortage or effectiveness/readiness level of end units that must be maintained. This in turn is related to (1) the level of repair and maintenance policy to be used by the contractor with existing equipments and support plans subject to deployment requirements; (2) the man loading constraints, quantity and quality; and (3) the transportation modes including priority rules and standardization times and costs.

That portion of the  $S_{ij}$  costs relating to follow-on costs reflect on both the amount of common assets and the operating decision rules used. Replenishment cost estimates as with policy rules must be based upon standardized operating management rules if support costs are to be comparable between competitor designs, and on the asset levels for common items which are required by competitor weapon systems. In effect, these decision rules, noted previously, control the replenishment of stocks in the pipeline which are, in turn, a function of not only demands on common resources among competitive systems, but also of policies on schedules maintenance, and overhaul, asset allocation among echelons, maintenance analysis and repair policies, and modes of transportation.†

### ILLUSTRATION OF THESIS

The support costing problem has been discussed in generic terms. An example might help to make the point. The reference is a U.S. Air Force Repair Level Decision Manual relating to the problem of "optimum" repair level decisions between intermediate and depot repair.

The "givens," to start with, consist of the following:

1. The item is unique — a landing gear strut assembly for a number of squadrons of a defined size. Expected and minimum monthly flying hours are given for a life cycle of 10 years. All aircraft are of like configurations and are delivered operational concurrently.

\*Implicit here is the question of the extent that contractors should be restricted in recommending equipment/systems to a specific scenario requirement and deployment plan for costing purposes. This holds for the entire design problem which ideally should be iterative between seller and buyer.

†Self-evident is the need for standardized costs for such resources as manpower, transportation, repair and storage facilities, and technical publications. Not so evident is the awareness that standardized costs can unwittingly favor one design over another--e.g., high standard costs for labor favor automated designs cost-wise. Cf. "The Resource Pricing Problem," N-464(R), June 1967, by W. Niskanen for a discussion on the result of errors in costing.



One-third of the reparable are assumed to be generated outside the United States. Cost of the assembly is given.

2. Repairs are made on the assembly with a given frequency based on aircraft usage.

3. Of these repairs, 70 percent can be done on the aircraft; of the remaining 30 percent pulled off, (a) 6 percent are condemned, (b) 12 percent will be repaired at depot, (c) 6 percent will be repaired at intermediate, and (d) 6 percent will be repaired either at intermediate or at depot.

The problem is for the contractor to recommend the repair level for the six percent in 3(d).

Additional information is provided in the manual to perform the calculations:

- The peculiar test equipment cost if repair is at depot (no maintenance cost is assumed for the 10-year life).

- Peculiar depot training; no additional costs for technical publications are assumed.

- If repair is at intermediate level, additional test equipment costs plus additional repair capability and maintenance costs for an existing test and repair station. Additional costs are also given for publications.

- Packaging and shipping costs for overseas and CONUS.

- Depot and intermediate turnaround time for repair for Hi-Value and regular items.

- Eight corrective tasks are defined with associated man-hours which is the basis for calculating expected labor costs at intermediate and depot levels.

- Consumable item cost for intermediate repair.

- Training costs for intermediate level.

From the above as a starting point (the "meeting of the minds" referred to at the outset of this paper) costs are calculated for depot versus intermediate repair for the pipeline test equipment transportation, safety level, training, and labor. But even these calculations require an understanding of rules to be followed, i.e., for safety level, spare parts based upon economic order quantity rules, turnaround times, and so forth.

In conclusion, of interest in this example is the fact that in order to make a choice of where to repair a unique item, the problem must be specified in terms of a large number of "givens," and even where calculations have to be made for operating costs, Air Force (DOD) policies and assumptions must be specified.

#### ACKNOWLEDGMENT

This paper was undertaken as a task under EIA G-426 Subcommittee on "Cost Methodology and Standardization." The following members of the Subcommittee reviewed and made contributions to this paper: Sgt. J. F. Clarke (USAF), D. Enfield (Lockheed), O. Gabrielson (Boeing), D. Gregor (Northrup), R. Highland (Hughes Aircraft), and J. Mosher (TRW).

\* \* \*

# A NOTE ON CYCLING IN THE SIMPLEX METHOD

K. T. Marshall\* and J. W. Suurballe

*Bell Telephone Laboratories Incorporated  
Holmdel, New Jersey*

## ABSTRACT

Although cycling in the simplex method has long been known, a number of theoretical questions concerning cycling have not been fully answered. One of these, stated in [3], is to find the smallest example of cycling, and Beale's example with three equations and seven variables is conjectured to be the smallest one. The exact bounds on dimensions of cycling examples are established in this paper. We show that Beale's example is the smallest one which cycles at a non-optimal solution, that a smaller one can cycle at the optimum, and that, in general (including the completely degenerate case), a cycling example must have at least two equations, at least six variables, and at least three non-basic variables. Examples and geometries are given for the extreme cases, showing that the bounds are sharp.

## 1. INTRODUCTION

It has long been known that cycling can occur in the Simplex method (see [1] and [2]). In [3] an unsolved problem is stated which is to find the smallest example of cycling. It is conjectured there that the example of Beale with three equations and seven variables is the smallest one.

We show that Beale's example is the smallest one which cycles at a point other than the optimum, but that a smaller one can cycle at the optimum. We prove that, in general (including the completely degenerate case), a cycling example must have at least six variables, at least two equations, and at least three nonbasic variables. Examples are given showing that these bounds are sharp, and geometries for the examples are included.

## 2. MAIN RESULT AND EXAMPLES

The following statements of the linear program and the pivot rule establish the setting for our analysis of cycling. Let the linear program be

$$\begin{aligned} & \text{minimize } c^T x^{(2)} \\ & \text{subject to } Ix^{(1)} + Ax^{(2)} = d \\ & x^{(1)}, x^{(2)} \geq 0, \end{aligned} \tag{1}$$

---

\*Presently in the Department of Operations Analysis, Naval Postgraduate School, Monterey, California.



where  $T$  denotes transpose,  $(I, A)$  is an  $(m \times n)$  matrix,  $x^{(1)}$  and  $d$  are  $m$ -vectors, and  $x^{(2)}$  and  $c$  are  $(n-m)$ -vectors. We assume that (1) is in canonical form: hence all the elements of  $d$  are non-negative, and  $I$  can be taken as the beginning basis for the primal Simplex method.

The pivot rule assumed is the standard one, with all possible ties resolved by the "leftmost" and "topmost" conventions:

(P1) the pivot column  $s$  is the leftmost column with

$$c_s = \min_j [c_j \mid c_j < 0];$$

(P2) the pivot row  $r$  is the topmost row with  $a_{rs} > 0$  and with

$$\frac{d_r}{a_{rs}} = \min_i \left[ \frac{d_i}{a_{is}} \mid a_{is} > 0 \right].$$

The possibility for cycling arises only in degenerate cases when  $d_r = 0$ .

We now state our main result.

**THEOREM:** Under the above rules in the primal Simplex method, for cycling to occur we must have  $m \geq 2$ ,  $n \geq m+3$ , and  $n \geq 6$ . For cycling to occur at a nonoptimal point we must have  $m \geq 3$ ,  $n \geq m+3$ , and  $n \geq 7$ . All the bounds are sharp.

We turn now to an informal discussion of these bounds and present examples of the extreme cases. Formal proofs of the results are given in Sections 3 and 4.

The bounds stated in the theorem are shown in Figure 1. They define a convex region in  $m, n$  space, with two extreme points  $(m, n) = (2, 6)$  and  $(3, 6)$  for cycling in general, and  $(m, n) = (3, 7)$  and  $(4, 7)$  for the case of cycling off the optimum.

**EXAMPLE 1:** The following two-equation problem in six variables illustrates one of the extreme points:

$$\begin{aligned} &\text{minimize } z = 2x_2 + 4x_4 + 4x_6 \\ &\text{subject to } x_1 \geq 0 \\ (2) \quad &x_1 - 3x_2 - x_3 - x_4 - x_5 + 6x_6 = 0 \\ &2x_2 + x_3 - 3x_4 - x_5 + 2x_6 = 0. \end{aligned}$$

This problem is put into canonical form with, say  $x_1$  and  $x_2$  as basic variables. The initial tableau in the Simplex method is then given by:

$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$-z$	
1		1/2	-11/2	-5/2	9		0
	1	1/2	-3/2	-1/2	1		0
		-1	7	1	2	1	0

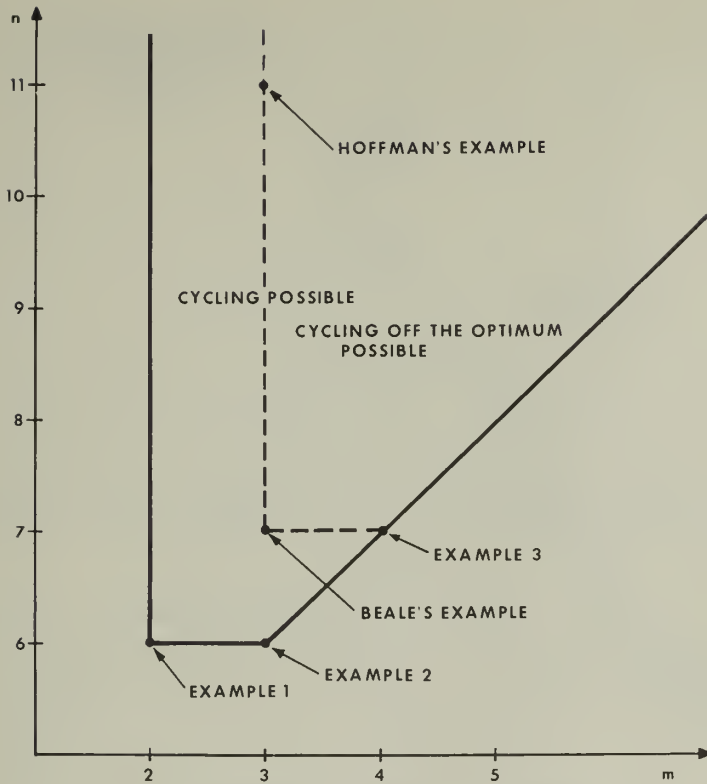


Figure 1. Bounds for cycling in the primal method

By following the above pivot rules it is seen that a sequence of six pivots will return the tableau to the above form. The six basic solutions are illustrated in Figure 2 in the space of the dual variables (the dual constraints are shown for the problem in its noncanonical form (2)). In this geometry each line represents a primal variable and the intersection of two lines represents a possible basic solution to the primal.\* The pivot rules choose the six basic solutions indicated in the figure, and these are seen to cycle around the optimal. It is interesting to note that the dual feasible region in this problem is bounded† (in fact it consists of a single point).

**EXAMPLE 2:** Next we give a three equation, six variable problem to show that the (3, 6) point is attainable. In canonical form the initial tableau is as follows:

\*Since the problem is completely degenerate all intersections correspond to basic feasible solutions.

†It can be shown in general for  $m=2$  that boundedness of the dual feasible region is necessary for cycling.

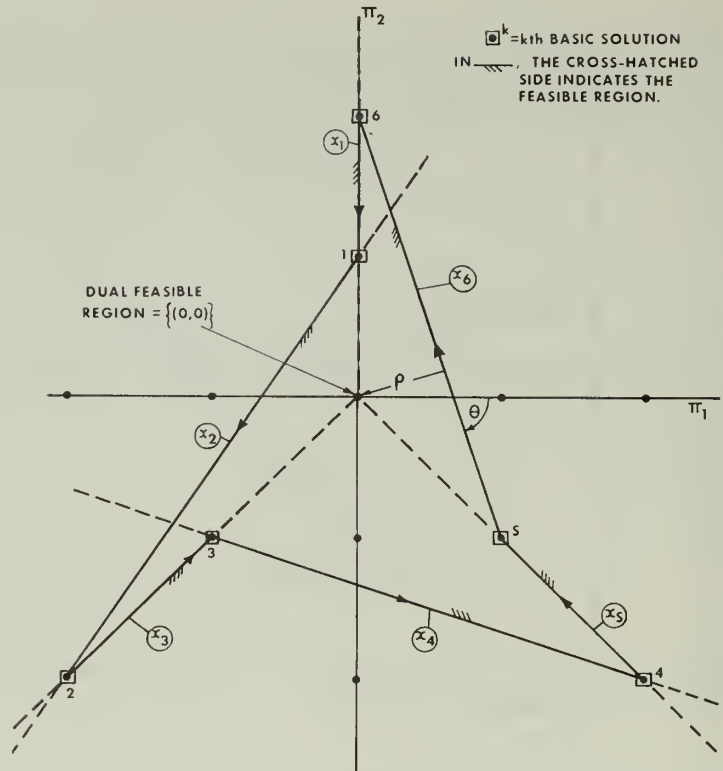


Figure 2. Cycling in Example 1

$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$-z$	
1			$3/5$	$-32/5$	$24/5$		0
	1		$1/5$	$-9/5$	$3/5$		0
		1	$2/5$	$-8/5$	$1/5$		0
			$-1/5$	$4/5$	$2/5$	1	0

By use of the above pivot rules the variables enter the basis in cyclic order. After three pivots the tableau is numerically the same, but with a cyclic permutation of columns. After six pivots the tableau returns to its original form. This example was actually derived from its dual. Its dual has an interesting geometrical interpretation which will be discussed later when we consider cycling in the dual Simplex method (see Example 4).

We now consider the extreme points of the region of problems which cycle off the optimum. Beale's example [2] is the case  $(m, n) = (3, 7)$ .

EXAMPLE 3: As an example of a  $(4, 7)$  problem which cycles off the optimum, we modify Example 2 as follows:

$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$	$-z$	
1				3/5	-32/5	24/5		0
	1			1/5	-9/5	3/5		0
		1		2/5	-8/5	1/5		0
			1		1			1
				-2/5	-2/5	9/5	1	0

Clearly, if  $x_6$  is brought into the basis a strict improvement in the objective function results and hence the problem cannot return to this basis. Under the above pivot rules, however, the sequence of basic solutions resemble those in Example 2 with  $x_4$  always basic.

### 3. BASIC LEMMAS AND GEOMETRY FOR THE CASE $m = 2$

In the theorem, the bounds  $m \geq 2$  and  $n \geq m+3$  imply the third bound  $n \geq 6$  for all cases except  $m = 2$ . Our approach will therefore be to establish the first two results for the general cycling case, then prove  $n \geq 6$  in the special case  $m = 2$ . This section contains lemmas and some special notation for the case  $m = 2$ .

In a completely degenerate two-equation problem (in canonical form)

(3)

$x_1$	$x_2$	$\dots$	$x_i$	$\dots$	$x_n$	$-z$	
1			$\dots$	$a$	$\dots$		0
			1	$\dots$	$b$	$\dots$	0
0	0	$\dots$	$c$	$\dots$		1	0

let  $a, b, c$  represent the row 1, row 2, and cost row elements, respectively, for an arbitrary

variable. For  $i = 1, 2$ , we use the notation  $\begin{pmatrix} a_i \\ b_i \\ c_i \end{pmatrix}$  and the name  $i$ -entry for a vector which has entered the basis through a row  $i$  pivot. Let

$$\bar{B} = \begin{bmatrix} a_1 & a_2 & 0 \\ b_1 & b_2 & 0 \\ c_1 & c_2 & 1 \end{bmatrix}$$

denote an arbitrary basis in this problem. Note that column order in  $\bar{B}$  may be different from the original column order in (3) for these basic vectors. Let  $B$  be the  $2 \times 2$  basis formed by deleting the last row and column in  $\bar{B}$ . We emphasize that non-primed letters always refer to the original coefficients in the above canonical form (3).

Let

$$\begin{pmatrix} a' \\ b' \\ c' \end{pmatrix} = \bar{B}^{-1} \begin{pmatrix} a \\ b \\ c \end{pmatrix}$$

represent the revised form for a vector in terms of the basis  $\bar{B}$ . In reference to the pivot rule (P2), a column  $\begin{pmatrix} a \\ b \\ c \end{pmatrix}$  can be the next 1-entry only if  $a' > 0$ , or can be the next 2-entry only if  $a' \leq 0, b' > 0$ .

The dual variables  $\pi_1$  and  $\pi_2$  are also the negative revised costs for  $x_1$  and  $x_2$ , respectively. A change in  $\pi_i$  for a given pivot operation is denoted by

$$\Delta\pi_i = \pi_i \text{ (new basis)} - \pi_i \text{ (old basis)}.$$

The following elementary relations are used:

$$(4) \quad \begin{cases} D = |B| = a_1 b_2 - a_2 b_1 \\ Da' = ab_2 - a_2 b \\ Db' = -ab_1 + a_1 b \end{cases}$$

$$(5) \quad \begin{cases} a = a_1 a' + a_2 b' \\ b = b_1 a' + b_2 b' \\ c = c' + c_1 a' + c_2 b' \end{cases}$$

$$(6) \quad \begin{cases} \Delta\pi_2 = -a_2 c' / a' D & \text{for a row 1 pivot} \\ \Delta\pi_1 = a_1 c' / b' D & \text{for a row 2 pivot} \end{cases}.$$

The lemmas now to be proved are used (in the proof of the Theorem) to characterize six distinct 1- and 2-entries in a cycle. They also give some geometric insight into the structure of cycling examples, at least for  $m = 2$ .



Let a column  $\begin{pmatrix} a \\ b \\ c \end{pmatrix}$  and its corresponding line in the dual space be described by a directed angle  $\theta = \arctan \frac{a}{b}$  and a directed distance to the origin  $\rho = \frac{c}{\sqrt{a^2 + b^2}}$ , as indicated

for variable  $x_6$  in Figure 2. With these concepts the lemmas establish that: in a sequence of pivots the lines corresponding to 2-entries are in strict counterclockwise progression, or are parallel with strict progress in the feasible direction. The example of Figure 2 is a case with strict counterclockwise progression for the 2-entries  $x_2, x_4, x_6$ , showing the complete 360-degree rotation which is necessary for 2-entries in a cycle.

**LEMMA 1:** In a given sequence of pivots for (3), suppose that  $b_2 \neq 0$  in the initial basis, and suppose that  $b_2$  does not change sign in the sequence of bases ( $b_2 = 0$  may occur). Then all 2-entries are distinct; in fact the sequence of vectors

$$\left( \frac{a_2}{b_2}, \frac{c_2}{\sqrt{a_2^2 + b_2^2}} \right)$$

derived from the 2-entries is strictly decreasing in the lexicographical sense (when  $b_2 = 0$ , we define  $\frac{a_2}{b_2} = -\infty$ ).

**LEMMA 2:** At a given basis  $\begin{pmatrix} a_1 & a_2 \\ b_1 & b_2 \end{pmatrix}$  for (3), suppose that the next entry  $\begin{pmatrix} a \\ b \end{pmatrix}$  is a 2-entry with  $bb_2 \leq 0$ . If  $b \neq 0$ , then  $b_1b < 0$ ; if  $b_2 \neq 0$ , then  $b_1b_2 > 0$ .

**LEMMA 3:** In a sequence of pivots for (3), the sign of the basis determinant  $D$  does not change.

**PROOF (for Lemma 3):** Using relations (4), with  $\begin{pmatrix} a \\ b \end{pmatrix}$  a new basis entry,  $D$  and  $D'$  the old and new basis determinants respectively, we get

$$\begin{aligned} D' &= Da' && \text{for a 1-pivot} \\ &= Db' && \text{for a 2-pivot.} \end{aligned}$$

By the standard pivot rule,  $a' > 0$  for a 1-pivot and  $b' > 0$  for a 2-pivot; therefore in all cases  $D$  and  $D'$  have the same sign.

**PROOF (for Lemma 2):** Because  $\begin{pmatrix} a \\ b \end{pmatrix}$  is a 2-entry, we have  $a' \leq 0$ ,  $b' > 0$ ; assume first that  $b \neq 0$  and multiply the second relation in (5) by  $b$  to get

$$bb_1a' = b^2 - bb_2b'.$$

From the various inequalities above, the right member is strictly positive, and since  $a' \leq 0$ , we conclude that  $a' < 0$  and  $b_1 b < 0$ . Assume next that  $b_2 \neq 0$  and multiply the second relation in (5) by  $b_2$  for a similar proof that  $b_1 b_2 > 0$ .

PROOF (for Lemma 1): Let the next basis entry  $\begin{pmatrix} a \\ b \\ c \end{pmatrix}$  be a 2-entry  $\Rightarrow [a' \leq 0, b' > 0]$

CASE 1 ( $a' < 0$ ): This implies  $ab_2 < a_2 b$  from (4). Also  $bb_2 > 0$ , or  $bb_2 = 0$  and  $b + b_2 \neq 0$ .

$$(i) \quad bb_2 > 0 \Rightarrow a_2/b_2 > a/b.$$

$$(ii) \quad b = 0 \Rightarrow b_2 \neq 0 \Rightarrow a_2/b_2 > -\infty = a/b.$$

$$(iii) \quad b_2 = 0 \text{ in some basis in the sequence with } a' < 0.$$

Then  $b \neq 0$  in the next basis  $\begin{pmatrix} a_1 & a \\ b_1 & b \end{pmatrix}$ , and there exists a last previous basis  $\begin{pmatrix} \tilde{a}_1 & \tilde{a}_2 \\ \tilde{b}_1 & \tilde{b}_2 \end{pmatrix}$  with  $\tilde{b}_2 \neq 0$ . We prove that  $\text{Sgn } b \neq \text{Sgn } \tilde{b}_2$ . This contradicts the hypothesis that no sign change in  $b_2$  occurs in the sequence.\* From Lemma 2,  $\tilde{b}_1 \tilde{b}_2 > 0$  and  $b_1 b < 0$ , so

$$(7) \quad \text{Sgn } \tilde{b}_2 = \text{Sgn } \tilde{b}_1 \text{ and } \text{Sgn } b \neq \text{Sgn } b_1.$$

The first basis with  $b_2 = 0$  has the form  $\begin{pmatrix} \tilde{a}_1 & a_2 \\ \tilde{b}_1 & 0 \end{pmatrix}$ , and the last basis with  $b_2 = 0$  has the form

$\begin{pmatrix} a_1 & a_2 \\ b_1 & 0 \end{pmatrix}$ . By Lemma 3,  $D = -a_2 b_1 > 0$  in all bases with  $b_2 = 0$ , so elements  $a_2$  and  $b_1$  are nonzero and have fixed signs for all bases in the sequence with  $b_2 = 0$ . In particular,  $\text{Sgn } \tilde{b}_1 = \text{Sgn } b_1$ . Hence from (7)  $\text{Sgn } \tilde{b}_2 \neq \text{Sgn } b$ .

CASE 2 ( $a' = 0$ ): From the second relation in (4),  $ab_2 = a_2 b$ , and since  $b$  and  $b_2$  do not have different signs, either  $bb_2 > 0$ , or  $b = b_2 = 0$ . If  $b \neq 0$ , then  $a/b = a_2/b_2$  and the ratio is finite, and if  $b = 0$ ,  $a/b = a_2/b_2 = -\infty$ . In any case,  $a_2/b_2$  remains constant when entries have  $a' = 0$ . We now show that

$$\frac{c_2}{\sqrt{a_2^2 + b_2^2}}$$

is strictly decreasing when entries have  $a' = 0$ . From (5),

\*Actually we show that in a sequence of bases with no sign change in  $b_2$ , if  $b_2 = 0$  ever occurs, all remaining bases in the sequence have  $b_2 = 0$ .

$$a = a_2 b'$$

$$b = b_2 b'$$

$$c = c' + c_2 b' < c_2 b' \quad \text{since } c' < 0.$$

Using these relations with the nonzero quantity  $\sqrt{a_2^2 + b_2^2}$ , we then have

$$c \sqrt{a_2^2 + b_2^2} < c_2 b' \sqrt{a_2^2 + b_2^2} = c_2 \sqrt{a^2 + b^2};$$

therefore,

$$\frac{c}{\sqrt{a^2 + b^2}} < \frac{c_2}{\sqrt{a_2^2 + b_2^2}}.$$

#### 4. PROOF OF MAIN THEOREM

It is easy to see that a single equation degenerate system cannot cycle. If  $x_i$  is the current basic variable, successive pivots will drive its cost coefficient positive and hence it will never again become a candidate for a pivot choice.

We now show that a system with only two nonbasic variables cannot cycle. Assume without loss of generality that the current tableau is of the following form, with  $c_{m+1} \leq c_{m+2}$  and  $c_{m+1} < 0$ ,  $d_i \geq 0$  for all  $i$ :

$x_1$	-	-	-	-	-	$x_m$	$x_{m+1}$	$x_{m+2}$	$-z$	
1	0	.....				0	$a_{11}$	$a_{12}$	0	$d_1$
	.						.	.	.	.
		.					.	.	.	.
			.				.	.	.	.
0	...	0	1	0	...	0	$a_{i1}$	$a_{i2}$	0	$d_i$
				.			.	.	.	.
					.		.	.	.	.
0	.....					0	$a_{m1}$	$a_{m2}$	0	$d_m$
0	-	-	-	-	-	0	$c_{m+1}$	$c_{m+2}$	1	

For cycling to occur we must find a finite sequence of pivots which return to this starting tableau. Obviously pivoting on any row with  $d_i > 0$  will reduce the value of the objective function and hence no cycling can occur. Since column  $(m+1)$  is the pivot column choose

$$r_1 = \min\{i \mid a_{i1} > 0, d_i = 0\}$$

and let  $S_1 = \{i \mid i \leq r_1\}$ . We pivot on  $a_{r_1 1}$  and obtain the tableau:

$x_1$	$\dots$	$x_{r_1}$	$\dots$	$x_m$	$x_{m+1}$	$x_{m+2}$	$-z$	
1	$\dots$	$t_1$	$\dots$	0	0	$a'_{12}$	0	$d_1$
		$\vdots$			$\vdots$	$\vdots$	$\vdots$	$\vdots$
0	$\dots$	$t_{r_1}$	$\dots$	0	1	$a'_{r_1 2}$	0	0
		$\vdots$			$\vdots$	$\vdots$	$\vdots$	$\vdots$
0	$\dots$	$t_m$	$\dots$	1	0	$a'_{m2}$	0	$d_m$
0	$\dots$	$c'_{r_1}$	$\dots$	0	0	$c'_{m+2}$	1	

If  $c'_{m+2} \geq 0$ , the algorithm ends. Hence, we assume  $c'_{m+2} < 0$  and find  $r_2 = \min \{i | a'_{i2} > 0 \text{ and } d_i = 0\}$ . Set  $S_2 = \{i | i \leq r_2\}$ . Now, if ...

(i)  $r_2 < r_1$ , then  $a_{r_2 1} > 0$  requires that  $r_1 \leq r_2$  by definition of the previous pivot row  $r_1$ . Hence  $a_{r_2 1} \leq 0$ ,  $\Rightarrow t_{r_2} \geq 0$ ,  $\Rightarrow$  pivoting on  $a'_{r_2 2}$  ends the algorithm.

(ii)  $r_2 = r_1$ , pivoting on  $a'_{r_2 2}$  ends the algorithm since  $t_{r_1} > 0$ .

(iii)  $r_2 > r_1$ , then it is possible that the modified cost coefficient under  $x_{r_1}$  goes negative and the algorithm continues.

Let  $|S_i|$  be the number of elements in set  $S_i$ , then we have shown that for the algorithm to continue  $|S_1| < |S_2|$ . Repeated application shows that  $|S_i| > |S_{i-1}|$  at each pivot  $i$ , and since  $1 \leq |S_1| \leq m$  the algorithm must end in at most  $m$  steps.

With the above proof that  $m \geq 2$  and  $n \geq m+3$ , the third bound  $n \geq 6$  is established for all cases except  $m = 2$ . For this case we assume that there is a cycle in the two-equation problem (3), and characterize six distinct columns  $x_1, \dots, x_6$  and six distinct bases  $B_1, \dots, B_6$ , which must be in the cycle. The logical construction of the pivot sequence can be followed in Table 1 below, where all the signs shown for the various elements are proved.

With (3) in canonical form, let  $x_1$  and  $x_2$  be the beginning 1- and 2-entries, respectively, and assume that a row 2 pivot occurs first. Then  $b_2 > 0$  in the first two bases, shown as  $B_1$  and  $B_2$  in Table 1. Since cycling is assumed to occur, there must eventually be a basis with  $b_2 < 0$ ; otherwise no sign change occurs in  $b_2$  and by Lemma 1 all the 2-entries are distinct. The first basis with  $b_2 < 0$  is indicated in Table 1 as  $B_4$ . Since the original basis eventually reappears, a new sequence with  $b_2 > 0$  must occur, and the first of these is indicated in Table 1 as  $B_6$ .

$B_4$  is defined above as the first basis with  $b_2 < 0$ , and we show by Lemma 2 that

$b_1 > 0$  in  $B_4$ : the last entry to  $B_4$  is a 2-entry  $\begin{pmatrix} a \\ b \end{pmatrix}$  with  $b < 0$ , the basis preceding  $B_4$  has  $b_2 \geq 0$ , so  $bb_2 \leq 0$  and all conditions of the lemma are satisfied. Therefore,  $bb_1 < 0$ , so  $b_1 > 0$ , with  $b_1$  common to  $B_4$  and the basis preceding it.

Similarly, the last entry to  $B_6$  is a 2-entry, and by Lemma 2 we have  $b_1 < 0$  in  $B_6$ .

TABLE 1  
Structure of Cycle for  $m = 2$

Basis Name	$B_1$	$B_2$	...	$B_3$	...	$B_4$	...	$B_5$	...	$B_6$
1-entries $\begin{pmatrix} a_1 \\ b_1 \end{pmatrix}$	$\begin{pmatrix} 1 \\ 0 \end{pmatrix} = x_1$	$\begin{pmatrix} 1 \\ 0 \end{pmatrix}$	...	$\begin{pmatrix} \cdot \\ + \end{pmatrix} = x_3$	...	$\begin{pmatrix} \cdot \\ + \end{pmatrix}$	...	$\begin{pmatrix} \cdot \\ - \end{pmatrix} = x_5$	...	$\begin{pmatrix} \cdot \\ - \end{pmatrix}$
2-entries $\begin{pmatrix} a_2 \\ b_2 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 1 \end{pmatrix} = x_2$	$\begin{pmatrix} \ominus \\ + \end{pmatrix}$	...	$\begin{pmatrix} \ominus \\ + \end{pmatrix} = x_4$	...	$\begin{pmatrix} \cdot \\ - \end{pmatrix}$	...	$\begin{pmatrix} \cdot \\ - \end{pmatrix} = x_6$	...	$\begin{pmatrix} \cdot \\ + \end{pmatrix}$
Type of pivot producing basis		row 2	...	row 1	...	row 2	...	row 1	...	row 2

The symbols  $+$ ,  $\ominus$ ,  $-$ ,  $\cdot$ , represent elements which are, respectively, positive, non-positive, negative, and unspecified.

$B_1$  = Initial basis

$B_2$  = Second basis

$B_3$  = First basis with  $b_1 > 0$

$B_4$  = First basis with  $b_2 < 0$

$B_5$  = First basis after  $B_4$  with  $b_1 < 0$

$B_6$  = First basis after  $B_4$  with  $b_2 > 0$

Other basis changes may occur in the sequence, as indicated between  $B_i$  and  $B_{i+1}$ ,  $i = 2, 3, 4, 5$ .

Since there exists a basis  $B_4$  with  $b_1 > 0$ , it must be preceded by a first basis with  $b_1 > 0$ , and it is called  $B_3$  in Table 1. Similarly, since  $B_6$  follows  $B_4$  and has  $b_1 < 0$ , there must be a first basis following  $B_4$  with  $b_1 < 0$ , and it is called  $B_5$  in Table 1. In Table 1 we let  $x_3$  and  $x_4$  be 1- and 2-entries in  $B_3$ , respectively, and  $x_5$  and  $x_6$  are 1- and 2-entries for  $B_5$ , respectively.

We now establish the signs for  $b_2$  in  $B_3$  and  $B_5$ : beginning at basis  $B_2$ , where  $b_2 > 0$ , Lemma 2 says that  $b_2$  must remain positive until after a basis with  $b_1 > 0$  occurs. Since  $B_3$  is the first basis with  $b_1 > 0$ , this means that  $b_2 > 0$  in  $B_3$  and in all preceding bases. Similarly,  $B_4$  has  $b_2 < 0$ , and we know from Lemma 2 that  $b_2$  must remain negative until after a basis with  $b_1 < 0$  occurs;  $B_5$  is the first basis (after  $B_4$ ) with  $b_1 < 0$ , so all bases from  $B_4$  to  $B_5$ , inclusive, have  $b_2 < 0$ .

We have shown that the columns  $x_1$  through  $x_6$  are necessarily in the cycle. Now we show that they are distinct. By Lemma 1, the 2-entries  $x_2$  and  $x_4$  are distinct because no sign change in  $b_2$  occurs between them. Although  $x_3$  and  $x_4$  both have  $b > 0$ , they form a basis and are independent. Similarly,  $x_5$  and  $x_6$  are independent. Dissimilarities in the elements  $b$  rule out all other possible equalities, except between  $x_2$  and  $x_3$ . We show  $x_2 \neq x_3$ , by showing that  $x_2$  and  $x_3$  have positive and negative revised costs, respectively, when  $x_3$  enters the basis. It is already known that  $b_2 > 0$  in the sequence of bases preceding  $B_3$ , so by Lemma 1,



$$\left( \frac{a_2}{b_2}, \frac{c_2}{\sqrt{a_2^2 + b_2^2}} \right) < (0,0)$$

in the lexicographical sense for each basis, and therefore  $a_2 \leq 0$  in the sequence.  $B_3$  is the first basis with  $b_1 > 0$ , so  $b_1 \leq 0$  in all earlier bases. Finally,  $D > 0$  is required by Lemma 3 because  $D > 0$  initially;  $a_1 > 0$  is the only solution. In summary,  $a_1 > 0$  and  $a_2 \leq 0$  in all bases up to  $B_3$ .

Let  $x_2$  have the revised cost  $-\pi_2$  and  $-\pi'_2$  for old and new bases, respectively. From (6), the new less old cost is

$$\begin{aligned} \pi_2 - \pi'_2 &= a_2 \frac{c'}{a'_D} && \text{for a row 1 pivot} \\ &= -a_1 \frac{c'}{b'_D} && \text{for a row 2 pivot.} \end{aligned}$$

Use of the signs found for  $a_1$  and  $a_2$  above means that the revised cost for  $x_2$  does not decrease for row 1 pivots and increases for row 2 pivots, up to the basis before  $B_3$ , so is strictly positive there. On the other hand  $x_3$  is the next entry, so it has a negative revised cost.

This completes the proof that in general,  $m \geq 2$ ,  $n \geq m+3$ ,  $n \geq 6$  is required for cycling in the primal Simplex method. We now prove the corresponding result  $m \geq 3$ ,  $n \geq m+3$ ,  $n \geq 7$  for cycling off the optimal solution.

First, there must be an element in the right-hand side which is not zero; otherwise all basic solutions have the value 0, and the cycle is at the optimum. To a nonzero element in the right-hand side there corresponds a nonzero basic variable, and a row containing it; this row and column cannot be pivot lines in the cycle, and thus every example cycling off the optimal value has at least one row and one column not in the cycle.

Sharpness of the bounds follows from Examples 1, 2, and 3, and from Beale's Example.

## 5. CYCLING IN THE DUAL SIMPLEX METHOD

We now give examples of the extreme cases for cycling in the dual Simplex method, which are duals to the extreme cases in the primal method. The barycentric coordinate geometry for the case dual to Example 2 in the primal method is shown.

The standard pivot rule for the dual method resolves ties by the "leftmost" and "rightmost" conventions, and is stated here using the generic terms  $d_{ij}$ ,  $\alpha_j$ ,  $\beta_i$  for the revised matrix, cost row, and right-hand side, respectively:

(D1) the pivot row  $r$  is the topmost row with

$$\beta_r = \min_i [\beta_i | \beta_i < 0]$$

(D2) the pivot column  $s$  is the leftmost column with  $d_{rs} < 0$  and with

$$\frac{\alpha_s}{(-d_{rs})} = \min_j \left[ \frac{\alpha_j}{(-d_{rj})} \mid d_{rj} < 0 \right]$$

A result can be stated for cycling in the dual method, similar to the theorem for cycling in the primal method; under the rules (D1) and (D2) for the dual simplex method, for cycling to occur we must have  $m \geq 3$ ,  $n \geq m+2$ , and  $n \geq 6$ . For cycling to occur at a nonoptimal value we must have  $m \geq 3$ ,  $n \geq m+3$ , and  $n \geq 7$ . All bounds are sharp.

An example which cycles in the primal method may also cycle in the dual method when transformed from the form of problem (1) into the problem

$$\begin{aligned} & \text{minimize } d^T y^{(2)} \\ (8) \quad & \text{subject to } Iy^{(1)} - A^T y^{(2)} = c \\ & y^{(1)}, y^{(2)} \geq 0, \end{aligned}$$

where elements of the vector  $y^{(2)}$  are the negative dual variables for problem (1). Here  $y^{(1)}$  and  $y^{(2)}$  are  $m$ -vectors, and  $(I, -A^T)$  is an  $(n-m) \times n$  matrix.  $I$  can be taken as a beginning basis for the dual simplex method since  $d$  is non-negative. We will call problems in forms (1) and (8) transforms of each other.

When stated in both forms (1) and (8), an example which cycles in one method (primal or dual) may or may not cycle in the other method when the standard pivot rule for each method is used. Examples of both types occur.

EXAMPLE 4: As mentioned previously, Example 2 is originally formulated from the dual point of view, and in this form has an interesting geometrical interpretation in barycentric coordinates. While Example 2 is an extreme case (3, 6) for cycling in the primal method, its original form given below is an extreme case (3, 6) for cycling in the dual method.

$$\begin{array}{llllll} \text{minimize} & y_1 & +y_2 & +y_3 & +y_4 & +y_5 & +y_6 \\ \text{subject to} & -3y_1 & & +y_3 & -3y_4 & & +y_6 = 1 \\ & 6y_1 & +y_2 & +y_3 & -2y_4 & & -y_6 = 0 \\ & -2y_1 & & -y_3 & +6y_4 & +y_5 & +y_6 = 0 \end{array}$$

and  $y_i \geq 0$ .

The beginning tableau given below with basic variables  $y_1, y_2, y_3$ , is the transform of the beginning tableau for Example 2.

$y_1$	$y_2$	$y_3$	$y_4$	$y_5$	$y_6$	$-z$	
1			$-3/5$	$-1/5$	$-2/5$		$-1/5$
	1		$32/5$	$9/5$	$8/5$		$4/5$
		1	$-24/5$	$-3/5$	$-1/5$		$2/5$
			0	0	0	1	0 (minimize $z$ )

By use of pivot rules (D1) and (D2) for the dual simplex method, the original tableau is realized in six pivots. Each column in this example has coefficients totaling 1, so columns can be plotted as points in the barycentric coordinate system of Figure 3. (Columns plotted are from the original form of the example above.) The unit total for column coefficients is preserved by pivots, so throughout all six pivots the revised tableaux give barycentric coordinates for points (columns) in terms of a basic triangle (current basis). In this geometry, a basic solution is optimal when the point corresponding to the right-hand side is in or on the basic triangle, and this never happens in the sequence of basic triangles 123, 234, ..., 561, 612.

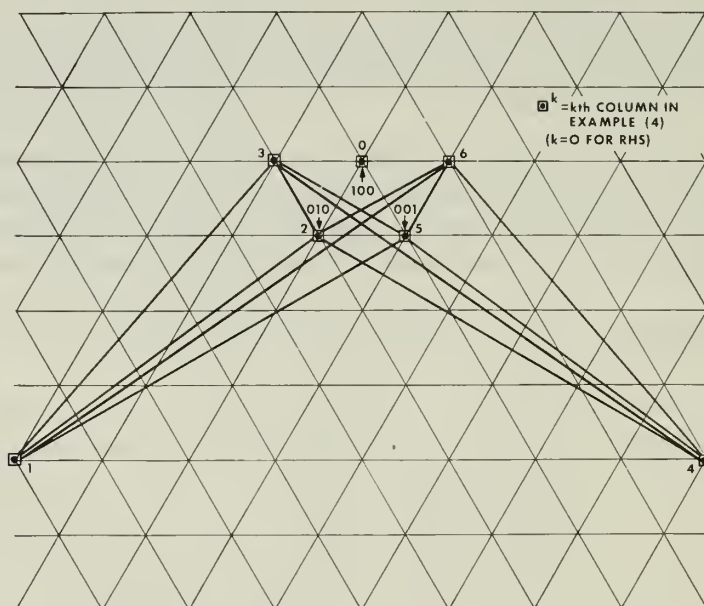


Figure 3. Cycling in Example 4: Barycentric coordinates with coordinate triangle  $\square^0, \square^2, \square^5$

EXAMPLE 5: This two-equation, six variable example in the form of (1)

$$\begin{aligned}
 &\text{minimize} && c_1 x_2 && + c_2 x_4 && + c_3 x_6 \\
 (9) \quad &\text{subject to} && x_1 && -\frac{1}{2} c_1 x_2 && - c_2 x_4 && - 2x_5 && + 2c_3 x_6 = 0 \\
 &&& c_1 x_2 && + x_3 && -\frac{3}{4} c_2 x_4 && - x_5 && + \frac{1}{2} c_3 x_6 = 0 \\
 &&& \text{and } x_j \geq 0
 \end{aligned}$$

is an extreme case for cycling in both primal and dual methods, and also illustrates cases in its range of parameters  $c_1, c_2, c_3$  where cycling need not occur in both methods. For cycling in the primal method, begin with the canonical form of (9) with  $x_1$  and  $x_2$  as basic variables. For cycling in the dual method, begin with the transform of this canonical form, as follows:

$y_1$	$y_2$	$y_3$	$y_4$	$y_5$	$y_6$	$-z$	
1				$-1/2$	$-1/c_1$		$-1$
	1			$(11/8)c_2$	$(3/4)(c_2/c_1)$		$(7/4)c_2$
		1		$5/2$	$1/c_1$		1
			1	$(-9/4)c_3$	$(-1/2)(c_3/c_1)$		$c_3/2$
				0	0	1	0

Cycling by either method depends on the parameters  $c_i$ , as indicated in (10) below.\* In the form of problem (1) this example cycles in six pivots in the primal method. In the dual method, six pivots in the above tableau return it to its original form except for a cyclic row permutation; six more pivots, for a total of twelve, are required to complete the cycle in the dual method. For cycling to occur, the costs  $c_1, c_2, c_3$  in Example 5 must lie in the following ranges:

$c_1$	$c_2$	$c_3$	Simplex Methods Which Cycle
$>2$	$>4$	$>2$	Both
2	4	$>2$	Primal, not dual
2	4	2	Neither

(10)

\*The conditions given in (10) are easily obtained from the pivot rules.

Examples of the two extreme points in cycling off the optimum are easily constructed from Examples 4 and 5 above.

## 6. ADDENDUM

Some results on cycling in [4] were brought to the attention of the authors by Saul I. Gass after most of this paper was finished. Essentially, Yudin and Gol'shtein show in [4] that for the primal simplex method to cycle the problem must have  $m \geq 2$ , and for  $m = 2$  there must be at least six iterations of the algorithm. They give necessary and sufficient conditions for a given  $2 \times 6$  array to cycle using more relaxed pivot rules than those given above. They show that a general class of problems which satisfy these conditions is given by:

(11)

$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$-z$	
1		$\beta$	$\alpha$	$\gamma$	$\delta$		0
	1	$\delta$	$\gamma$	$\alpha$	$\beta$		0
		$-a$	$a$	$-a$	$a$	1	

where  $a > 0$ ,  $\delta > -\gamma > -\alpha > \beta > 0$ .

It is easily demonstrated, however, that such a problem has a nonfeasible dual (unbounded primal) and can be made into a bounded problem only by the addition of a seventh variable and third equation. Hence the question of whether or not a  $2 \times 6$  example of cycling could be found with a bounded solution was left unanswered.

It is easy to check that our Example 1 satisfies all the necessary and sufficient conditions in [4] (page 245). In attempting to modify their class of problems (11) to obtain a class with bounded solutions, we discovered the following class which satisfies all the necessary and sufficient conditions in [4], has bounded primal solutions, has unique dual solutions, and each member of the class terminates with an optimal solution after 2 pivots of the primal simplex algorithm under our pivot rules:

$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$-z$	
1		$\beta$	$\alpha$	$(\alpha - \gamma - \delta)$	$\delta$		0
	1	$\delta$	$\gamma$	$-\delta$	$\beta$		0
		$-a$	$a$	$a$	$a$	1	

where  $a > 0$ ,  $\delta > -\gamma > -\alpha > \beta > 0$ .



## REFERENCES

- [1] Hoffman, A. J., "Cycling in the Simplex Algorithm," National Bureau of Standards Report No. 2974 (Dec. 1953).
- [2] Beale, E. M. L., "Cycling in the Dual Simplex Algorithm," Nav. Res. Logist. Quart., 2, 269-275 (Dec. 1955).
- [3] Dantzig, G. B., Linear Programming and Extensions (Princeton University Press, Princeton, New Jersey, 1963), pp. 228-239.
- [4] Yudin, D. B. and E. G. Gol'shtein, Linear Programming (Israel Program of Scientific Translations, Jerusalem, 1965), pp. 238-250.

\*       \*       \*



## COMMUNICATION

### A NOTE ON A MAXIMUM UTILITY SOLUTION TO A VEHICLE CONSTRAINED TANKER SCHEDULING PROBLEM

M. Bellmore  
*The Johns Hopkins University*

G. Bennington  
*The Mitre Corporation*

S. Lubore  
*The Mitre Corporation*

Due to an oversight in our paper [1], we did not reference the work of Holladay [3] who also solved a vehicle constrained tanker scheduling problem. His matrix marking algorithm uses more computer storage than the method of our paper if the deliveries to be made consist of more than one vehicle load. The approach presented in our paper formalizes the realization of the scheduling problem as a time expanded network, which allows one to use standard algorithms such as described by Ford and Fulkerson [2]. The authors apologize for the oversight.

#### REFERENCES

- [1] Bellmore, M. G. Bennington, and S. Lubore, "A Maximum Utility Solution to a Vehicle Constrained Scheduling Problem," *Naval Research Logistics Quarterly*, 15, 403-411 (1968).
- [2] Ford, L. R. and D. R. Fulkerson, Flows in Networks (Princeton University Press, Princeton, N. J., 1962).
- [3] Holladay, J., "Some Transportation Problems and Techniques For Solving Them," *Naval Research Logistics Quarterly*, 11, 15-42 (1964).

\* \* \*



NATO RELIABILITY CONFERENCE

As part of its program for 1969 the NATO Advisory Panel on Operational Research is organizing, under the aegis of the NATO Science Committee, a Conference on the Application of Operational Research to Reliability. The Conference will be held in Turin, Italy from June 30 to July 4, 1969.

The theme of the conference is very broad. Reliability has already been the subject of many congresses in all countries.

This conference will deal more precisely with the operational research aspects of reliability.

- (a) Cost and effectiveness of the different methods for evaluating, controlling and improving the reliability of equipment.
- (b) Information systems: collection and analysis of the data concerning the life and failures of materials.
- (c) The relation between the reliability of a system and the reliability of its constituents: functions of structures, simulation of systems, redundancy.
- (d) Reliability prediction of a material before its manufacture. The methods of estimation will depend on the information already existing on the material: broad definition of the project, detailed plans, prototypes.
- (e) Reliability tests. The testing of a material in its normal working conditions requires often a long period of time so as to give a satisfactory estimation of its reliability; sometimes accelerated tests are given, under very hard working conditions.
- (f) Stocks of component parts. The stocks of component parts bought with the materials and delivered to the maintenance services must be established according to the reliability of the different constituents.
- (g) Preventive maintenance policy. The frequency of the tests and the preventive replacement must be calculated according to the wear of the materials (life distribution) and must be revised according to the operational results.
- (h) Modifications: The deficiencies shown by the tests or the operational functioning of the materials can often be eliminated by technical changes. The first changes to be made are those for the improvement of reliability at a lesser cost. Then, it is necessary to check the efficiency of these changes.

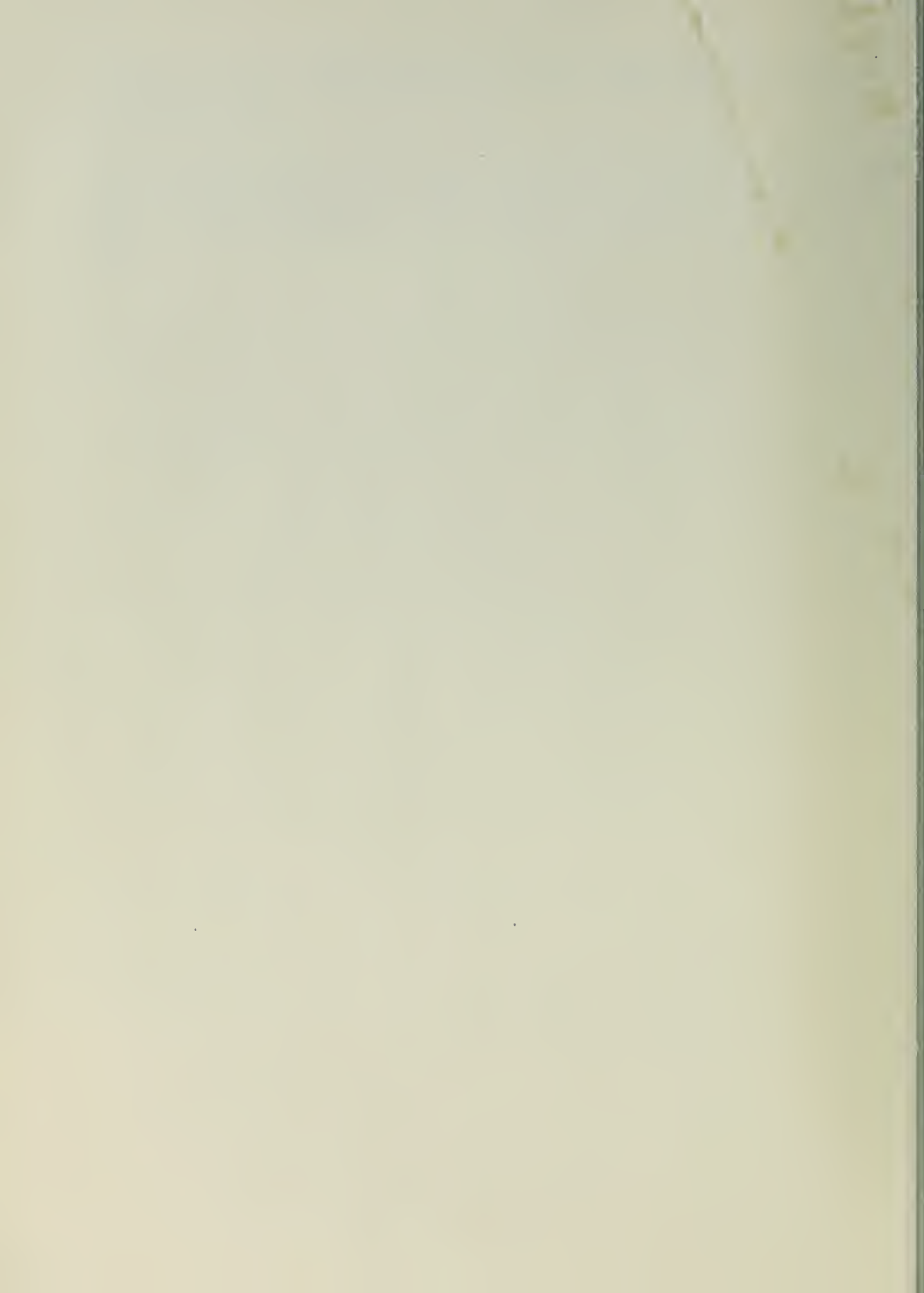


Among the operational research techniques most commonly used in reliability we will cite: statistical control tests, sampling plans, stochastic processes, simulation, Boolean algebra, linear programming.

The American point of contact for the conference is Dr. Ernest M. Scheuer, The Rand Corporation, 1700 Main Street, Santa Monica, California 90406. May 1, 1969 is the deadline for registration for Conference participants. Further details may be obtained from Dr. Scheuer at Rand.

\* \* \*





## INFORMATION FOR CONTRIBUTORS

The NAVAL RESEARCH LOGISTICS QUARTERLY is devoted to the dissemination of scientific information in logistics and will publish research and expository papers, including those in certain areas of mathematics, statistics, and economics, relevant to the over-all effort to improve the efficiency and effectiveness of logistics operations.

Manuscripts and other items for publication should be sent to The Managing Editor, NAVAL RESEARCH LOGISTICS QUARTERLY, Office of Naval Research, Washington, D.C. 20360. Each manuscript which is considered to be suitable material for the QUARTERLY is sent to one or more referees.

Manuscripts submitted for publication should be typewritten, double-spaced, and the author should retain a copy. Refereeing may be expedited if an extra copy of the manuscript is submitted with the original.

A short abstract (not over 400 words) should accompany each manuscript. This will appear at the head of the published paper in the QUARTERLY.

There is no authorization for compensation to authors for papers which have been accepted for publication. Authors will receive 50 reprints of their published papers.

Readers are invited to submit to the Managing Editor items of general interest in the field of logistics, for possible publication in the NEWS AND MEMORANDA or NOTES sections of the QUARTERLY.

---

---

CONTENTS

ARTICLES	Page
On a New Approach to the Analysis of Stationary Inventory Problems Oscar A. Z. Leneman and Frederick J. Beutler	1
The Status and Impact of Reliability Methodology Gerald J. Lieberman	17
A New Derivation of the Logistic Distribution Satya D. Dubey	37
On the Theory of Semi-Infinite Programming and a Generalization of the Kuhn-Tucker Saddle Point Theorem for Arbitrary Convex Functions A. Charnes, W. W. Cooper, and K. O. Kortanek	41
On Two Nonprobabilistic Utility Measures for Weapon Systems Hermann Enzer	53
An Evaluation of Incentive Contracting Experience I. N. Fisher	63
Inventory Control of By-Products Richard V. Evans	85
Discounted Production Scheduling and Employment Smoothing Steven A. Lippman and John S. C. Yuan	93
Problems in Life Cycle Support Cost Estimation A. S. Goldman	111
A Note on Cycling in the Simplex Method K. T. Marshall and J. W. Suurballe	121
COMMUNICATION	139
NEWS AND MEMORANDA	141

---

---